

Arbitrary-order mixed methods for heterogeneous anisotropic diffusion on general meshes

Daniele A. Di Pietro^{*1} and Alexandre Ern^{†2}

¹University of Montpellier, Institut Montpellierain Alexander Grothendieck, 34095 Montpellier, France

²Université Paris-Est, CERMICS (ENPC), 6–8 avenue Blaise Pascal, 77455, Marne la Vallée cedex 2, France

August 1, 2015

Abstract

We devise mixed methods for heterogeneous anisotropic diffusion problems supporting general polyhedral meshes. For a polynomial degree $k \geq 0$, we use as potential degrees of freedom the polynomials of degree at most k inside each mesh cell, whereas for the flux we use both polynomials of degree at most k for the normal component on each face and fluxes of polynomials of degree at most k inside each cell. The method relies on three ideas: a flux reconstruction obtained by solving independent local problems inside each mesh cell, a discrete divergence operator with a suitable commuting property, and a stabilization enjoying the same approximation properties as the flux reconstruction. Two static condensation strategies are proposed to reduce the size of the global problem, and links to existing methods are discussed. We carry out a full convergence analysis yielding flux-error estimates of order $(k + 1)$ and L^2 -potential estimates of order $(k + 2)$ if elliptic regularity holds. Numerical examples confirm the theoretical results.

1 Introduction

Over the last few years, several discretization methods for elliptic PDEs on general meshes including polyhedral cells and nonmatching interfaces have been proposed and analyzed. Such general meshes are useful, for instance, in the context of subsurface flow simulations in saline aquifers and oil reservoirs featuring geological layers with complex three-dimensional shapes. Another motivation for using general meshes stems from agglomeration-based mesh coarsening strategies [1, 4]. Examples of low-order discretization methods supporting general meshes include Mimetic Finite Differences (MFD) [40, 18, 15], Mixed Finite Volumes (MFV) [34] and Hybrid Finite Volumes (HFV) [37], the generalized Crouzeix–Raviart method [32], Gradient schemes [36], Cell-Centered Galerkin (CCG) methods [27], the Discrete Geometric Approach [26], and Compatible Discrete Operator (CDO) schemes [12, 13]. Tight connections exist between various of the above methods, as discussed in [35, 38, 33, 12]. Higher-order

^{*}daniele.di-pietro@umontpellier.fr, corresponding author

[†]ern@cermics.enpc.fr

discretizations extending ideas from the above methods have become recently available, and include high-order MFD [8, 9, 41], the Virtual Element (VE) method [5, 6], Weak Galerkin (WG) schemes [43, 44], and Hybrid High-Order (HHO) methods [30, 29, 31].

In many applications involving elliptic PDEs, in particular with heterogeneous and possibly anisotropic diffusion, it is often of interest to approximate both the potential (primal variable) and the diffusive flux (dual variable) starting from the elliptic PDE in mixed form. An extensive choice of mixed finite elements is available on simplicial and rectangular meshes, see [10] and references therein, while extensions to pyramids and triangular prisms can be found in [14]. The literature is more scarce on meshes with more general cell shapes. One possibility based on standard finite element spaces is to introduce a simplicial submesh of the polyhedral cells so as to reconstruct the diffusive flux from its normal interface values inside the cell by solving a local minimization problem [21, 42]. In the present work, we introduce and analyze an alternative and simpler mixed method on general meshes, avoiding the need for local solves on the submesh.

The key ideas can be summarized as follows. Let $k \geq 0$ be an arbitrary polynomial degree. As a starting point, we consider potential degrees of freedom (DoFs) that are polynomials of degree at most k in each cell, while the flux DoFs consist of face-based DoFs that are scalar-valued polynomials of degree at most k in each face (approximating the normal flux across the face) and of cell-based DoFs that are fluxes associated with potential polynomials of degree at most k . Then, we devise two key discrete operators, both acting cell-wise: (i) a discrete divergence operator, mapping flux DoFs to potential DoFs, that satisfies a commuting property with suitable reduction operators acting on continuous fluxes and potentials, and (ii) a consistent flux reconstruction operator mapping flux DoFs to a continuous flux evaluated from the gradient of a polynomial potential of degree at most $(k + 1)$; consistency stems from the fact that the flux reduction operator is a right inverse of the reconstruction operator. The third key ingredient is a stabilization bilinear form which is defined cell-wise as a least-squares penalty on cell faces between flux face-based DoFs and the normal component of the reconstructed flux. At the discrete level, before the static condensation discussed below, the weighted L^2 -inner product between discrete fluxes can be interpreted as originating from a discrete Hodge inner product built from a consistent and a stabilization part in the spirit of [19, 16, 37, 11]. One salient difference is that, in the present method, the reconstructed flux is not made fully explicit (only the consistent part is, but not the stabilization part).

A related approach that appeared recently is the mixed VE method of [17]. Therein, the authors propose a DoF space for the flux variable that allows one to compute the L^2 -projection on the space of polynomials of degree $\leq k$ of vector-valued fields from their DoFs, see also [7] for a different definition of some DoFs facilitating the computation of the L^2 -projection. One relevant difference is that the present construction uses fewer cell-based flux DoFs, the reduction factor being essentially equal to the space dimension, since these DoFs are evaluated herein as potential gradients and not as full vector fields. The second difference is that the consistent part of the formulation is devised as a flux reconstruction evaluated from a local linear system and is therefore fully computable instead of being virtual. This also leads to a different viewpoint in stabilization design. Moreover, the lowest-order case ($k = 0$) can be incorporated in a straightforward manner.

In practice, further computational savings can be achieved by static condensation, whereby cell-based DoFs are eliminated leading to a global linear system in terms of face-based DoFs. Two strategies for static condensation can actually be considered. In the first strategy, cell-based flux DoFs and potential DoFs up to a constant value per cell are eliminated locally.

The global problem is then of saddle-point form and involves the face-based flux DoFs and the mean-value of the potential in each cell. This problem has the same size and structure as that derived in the Multiscale Hybrid-Mixed (MHM) method [2, 39] using a different viewpoint based on the primal mixed method with Lagrange multipliers enforcing the interface continuity of the potential. In the second, more computationally-effective strategy, the face-based flux DoFs can be hybridized by means of face-based polynomials of degree at most k which play the role of Lagrange multipliers and can be interpreted as potential traces on cell faces. In this case, the flux and potential DoFs can be eliminated locally, and the global linear system in the Lagrange multipliers is symmetric positive definite. The size and structure of this system are closely related to those obtained using HHO methods [30, 29, 31] and also Hybridizable Discontinuous Galerkin (HDG) methods [23, 24, 25]. One difference with HDG is that the cell-based flux DoFs to be eliminated locally are not vector-valued polynomials but gradients of scalar-valued polynomials as in HHO methods, leading to a reduction factor in the number of local DoFs essentially equal to the space dimension, as mentioned above. We also emphasize that the present stabilization design is novel, and that bridges with the HDG methods exploiting this novel design are discussed in [22].

Interestingly, on matching simplicial (or Cartesian) meshes, the present method has also fewer local flux DoFs than classical (Raviart–Thomas or Brezzi–Douglas–Marini) mixed finite elements. This comes at the (affordable) price of solving one small linear problem per element to compute the consistent part of the flux reconstruction for polynomial degrees $k \geq 1$, while the consistent part can be computed explicitly in the lowest-order case $k = 0$.

This paper is organized as follows. In Section 2, we specify the continuous and discrete settings including the key ingredients to formulate the discrete problem. In Section 3, we present the discrete problem and state our main results which include flux-error estimates of order $(k + 1)$ and L^2 -potential estimates of order $(k + 2)$ if elliptic regularity holds. We also discuss static condensation, which is important in the practical implementation of the method, as well as links with other methods in the lowest-order case ($k = 0$). In particular, up to an equivalent design of the stabilization, the lowest-order Raviart–Thomas method is recovered if the mesh is made of simplices, and a similar observation can be made on Cartesian meshes. Interestingly, this link shows that Raviart–Thomas basis functions can be decomposed into a consistent and a stabilization part. In Section 4, we collect the proofs of our results. Finally, in Section 5, we present numerical results illustrating the error analysis.

2 Continuous and discrete settings

In this section, we introduce the continuous and discrete settings. In particular, we define the flux and potential DoFs, the discrete divergence and flux reconstruction operators, and the discrete Hodge inner product.

2.1 Continuous setting

Let $\Omega \subset \mathbb{R}^d$, $d \geq 2$, be an open, bounded, connected set with polygonal (or polyhedral) boundary. We consider the diffusion problem

$$\begin{aligned} -\nabla \cdot (\mathbf{K} \nabla u) &= f && \text{in } \Omega \\ u &= 0 && \text{on } \partial\Omega, \end{aligned} \tag{1}$$

where we enforce a homogeneous Dirichlet boundary condition for simplicity. The source term f is in $L^2(\Omega)$, and the diffusion tensor \mathbf{K} is piecewise constant and takes symmetric positive definite values with eigenvalues in the interval $[\lambda_b, \lambda_\sharp]$ with $0 < \lambda_b \leq \lambda_\sharp < +\infty$. For $X \subset \Omega$, we denote by $(\cdot, \cdot)_X$ and $\|\cdot\|_X$ respectively the standard inner product and norm of $L^2(X)$, with the convention that the index is omitted if $X = \Omega$. The same notation is used for vector-valued functions. Letting $\Sigma := \mathbf{H}(\text{div}; \Omega)$ and $U := L^2(\Omega)$, the mixed variational formulation of problem (1) reads: Find $(\boldsymbol{\sigma}, u) \in \Sigma \times U$ such that

$$(\mathbf{K}^{-1}\boldsymbol{\sigma}, \boldsymbol{\tau}) + (u, \nabla \cdot \boldsymbol{\tau}) = 0 \quad \forall \boldsymbol{\tau} \in \Sigma, \quad (2a)$$

$$(\nabla \cdot \boldsymbol{\sigma}, v) = -(f, v) \quad \forall v \in U. \quad (2b)$$

Throughout this work, $\boldsymbol{\sigma}$ and u are termed flux and potential, respectively. Boldface fonts are used for vector- and tensor-valued quantities.

2.2 Meshes and analysis tools

Denote by $\mathcal{H} \subset \mathbb{R}_*^+$ a countable set of mesh sizes having 0 as its unique accumulation point. Following [28, Chapter 1], we consider h -refined mesh sequences $(\mathcal{T}_h)_{h \in \mathcal{H}}$ where, for all $h \in \mathcal{H}$, \mathcal{T}_h is a finite collection of nonempty disjoint open polyhedral cells T such that $\bar{\Omega} = \bigcup_{T \in \mathcal{T}_h} \bar{T}$ and $h = \max_{T \in \mathcal{T}_h} h_T$ with h_T standing for the diameter of the cell T . Our analysis hinges on the following assumption on the mesh sequence.

Assumption 1 (Admissible mesh sequence). *For all $h \in \mathcal{H}$, \mathcal{T}_h admits a matching simplicial submesh \mathfrak{T}_h such that any cell and any face in \mathfrak{T}_h belongs to only one cell and face of \mathcal{T}_h , and there exists a real number $\varrho > 0$ independent of h such that, for all $h \in \mathcal{H}$, (i) for all simplex $S \in \mathfrak{T}_h$ of diameter h_S and inradius r_S , $\varrho h_S \leq r_S$ and (ii) for all $T \in \mathcal{T}_h$, and all $S \in \mathfrak{T}_h$ such that $S \subset T$, $\varrho h_T \leq h_S$.*

The simplicial submesh in Assumption 1 is just an analysis tool, and it is not used in the actual construction of the discretization method. Furthermore, all the meshes in \mathcal{T}_h are assumed to be compatible with the known partition on which the diffusion tensor is piecewise constant. The (constant) restriction of \mathbf{K} to a mesh cell $T \in \mathcal{T}_h$ is denoted \mathbf{K}_T . The lowest and largest eigenvalue of \mathbf{K}_T are denoted $\lambda_{b,T}$ and $\lambda_{\sharp,T}$, respectively, and we introduce the local anisotropy ratio $\rho_{\mathbf{K},T} := \lambda_{\sharp,T}/\lambda_{b,T}$. More generally, a piecewise-smooth diffusion tensor can be considered, but this would entail additional technicalities; in particular, Lemma 3 below on polynomial preservation no longer holds exactly, but up to a high-order term depending on the local Lipschitz constant of \mathbf{K}_T , similarly to the setting in [30] for HHO methods.

A face F is defined as a hyperplanar closed connected subset of $\bar{\Omega}$ with positive $(d-1)$ -dimensional Hausdorff measure and such that (i) either there exist $T_1, T_2 \in \mathcal{T}_h$ such that $F \subset \partial T_1 \cap \partial T_2$ and F is called an interface or (ii) there exists $T \in \mathcal{T}_h$ such that $F \subset \partial T \cap \partial \Omega$ and F is called a boundary face. Interfaces are collected in the set \mathcal{F}_h^i , boundary faces in \mathcal{F}_h^b , and we let $\mathcal{F}_h := \mathcal{F}_h^i \cup \mathcal{F}_h^b$. The diameter of a face $F \in \mathcal{F}_h$ is denoted by h_F . For all $T \in \mathcal{T}_h$, $\mathcal{F}_T := \{F \in \mathcal{F}_h \mid F \subset \partial T\}$ denotes the set of faces contained in ∂T (with ∂T denoting the boundary of T) and, for all $F \in \mathcal{F}_T$, \mathbf{n}_{TF} is the unit normal to F pointing out of T . For each interface $F \in \mathcal{F}_h^i$, we fix once and for all the ordering for the cells $T_1, T_2 \in \mathcal{T}_h$ such that $F \subset \partial T_1 \cap \partial T_2$ and we let $\mathbf{n}_F := \mathbf{n}_{T_1,F}$. For a boundary face, we simply take $\mathbf{n}_F = \mathbf{n}$, the outward unit normal to Ω . In what follows, $|\cdot|_l$ denotes the l -dimensional Hausdorff measure.

We recall some results that hold uniformly in h on admissible mesh sequences [28, §1.4]. For all $h \in \mathcal{H}$, all $T \in \mathcal{T}_h$, and all $F \in \mathcal{F}_T$, h_F is comparable to h_T :

$$\varrho^2 h_T \leq h_F \leq h_T. \quad (3)$$

Moreover, there exists an integer N_ϱ depending on ϱ such that

$$\forall h \in \mathcal{H}, \quad \max_{T \in \mathcal{T}_h} \text{card}(\mathcal{F}_T) \leq N_\varrho. \quad (4)$$

Let $l \geq 0$ be a non-negative integer. For an n -dimensional subset X of $\bar{\Omega}$ ($n \leq d$), we introduce the space $\mathbb{P}_n^l(X)$ spanned by the restriction to X of n -variate polynomials of total degree $\leq l$. Then, there exists a real number C_{tr} depending on ϱ and l , but independent of h , such that the following discrete trace inequality holds for all $T \in \mathcal{T}_h$ and all $F \in \mathcal{F}_T$:

$$\|v\|_F \leq C_{\text{tr}} h_F^{-1/2} \|v\|_T \quad \forall v \in \mathbb{P}_d^l(T). \quad (5)$$

Furthermore, the following inverse inequality holds for all $T \in \mathcal{T}_h$ with C_{inv} again depending on ϱ and l , but independent of h :

$$\|\nabla v\|_T \leq C_{\text{inv}} h_T^{-1} \|v\|_T \quad \forall v \in \mathbb{P}_d^l(T). \quad (6)$$

Moreover, there exists a real number C_{app} depending on ϱ and l , but independent of h , such that, for all $T \in \mathcal{T}_h$, denoting by π_T^l the L^2 -orthogonal projector on $\mathbb{P}_d^l(T)$, the following holds: For all $s \in \{1, \dots, l+1\}$ and all $v \in H^s(T)$,

$$|v - \pi_T^l v|_{H^m(T)} + h_T^{1/2} |v - \pi_T^l v|_{H^m(\partial T)} \leq C_{\text{app}} h_T^{s-m} |v|_{H^s(T)} \quad \forall m \in \{0, \dots, s-1\}. \quad (7)$$

Finally, the following Poincaré inequality is valid for all $T \in \mathcal{T}_h$ and all $v \in H^1(T)$ such that $\int_T v = 0$:

$$\|v\|_T \leq C_P h_T \|\nabla v\|_T, \quad (8)$$

where $C_P = \pi^{-1}$ for convex cells while, for more general cell shapes, C_P can be estimated in terms of ϱ .

In what follows, the regularity assumptions in the error estimates are expressed in terms of the broken Sobolev spaces $H^l(\mathcal{T}_h) := \{v \in L^2(\Omega) \mid v|_T \in H^l(T), \forall T \in \mathcal{T}_h\}$. Additionally, we often abbreviate as $A \lesssim B$ the inequality $A \leq cB$ with generic constant c uniform with respect to the mesh size and the diffusion tensor; the constant c can depend on the polynomial degree k .

Remark 1 (Face degeneration). *The present setting for mesh regularity requires that mesh faces have a comparable diameter to that of the cells they belong to, as reflected by the bounds in (3). This setting allows us to use the discrete trace inequality (5) and to use the length scale h_F in the design of the stabilization bilinear form, see (16) below. A framework allowing for face degeneration (keeping the bound (4)) has been proposed in [20] in the context of interior-penalty discontinuous Galerkin methods, allowing one to use a sharper discrete trace inequality on faces belonging to mesh cells matching the assumptions stated in [20, Def. 4.3]. In principle, one could expect that this framework could be used herein while adapting accordingly the penalty strategy; a careful inspection of this point is postponed to future work.*

2.3 Local degrees of freedom and reduction operator

Let $k \geq 0$. On every cell $T \in \mathcal{T}_h$, the DoFs for the flux and the potential are

$$\underline{\Sigma}_T^k := \mathbf{\Gamma}_T^k \times \left\{ \bigotimes_{F \in \mathcal{F}_T} \mathbb{P}_{d-1}^k(F) \right\}, \quad U_T^k := \mathbb{P}_d^k(T), \quad (9)$$

with $\mathbf{\Gamma}_T^k := \mathbf{K}_T \nabla \mathbb{P}_d^k(T)$. A generic collection of DoFs in $\underline{\Sigma}_T^k$ is denoted $\underline{\tau}_T = (\tau_T, (\tau_F)_{F \in \mathcal{F}_T})$. For $k = 0$, only the face-based flux DoFs are relevant.

Set $\Sigma^+(T) := \{\tau \in \mathbf{L}^s(T) \mid \nabla \cdot \tau \in L^2(T)\}$ with $s > 2$. The reduction operator $\underline{I}_T^k : \Sigma^+(T) \rightarrow \underline{\Sigma}_T^k$ is such that, for all $\tau \in \Sigma^+(T)$, $(\underline{I}_T^k \tau)_T = \mathbf{K}_T \nabla v$ where $v \in \mathbb{P}_d^k(T)$ solves the following Neumann problem:

$$((\underline{I}_T^k \tau)_T, \nabla w)_T = (\mathbf{K}_T \nabla v, \nabla w)_T = (\tau, \nabla w)_T \quad \forall w \in \mathbb{P}_d^k(T), \quad (10)$$

while $(\underline{I}_T^k \tau)_F = \pi_F^k(\tau \cdot \mathbf{n}_F)$ for all $F \in \mathcal{F}_T$, where π_F^k is the standard L^2 -orthogonal projector onto $\mathbb{P}_{d-1}^k(F)$. The Neumann problem (10) has compatible right-hand side vanishing for constant w , and its solution v is defined up to a constant, which we can fix by prescribing its average value on T . Additionally, the definition of $\Sigma^+(T)$ ensures that $(\underline{I}_T^k \tau)_F$ is well-defined.

2.4 Discrete divergence

The discrete divergence operator $D_T^k : \underline{\Sigma}_T^k \rightarrow U_T^k$ is such that, for all $(\underline{\tau}_T, v) \in \underline{\Sigma}_T^k \times U_T^k$,

$$(D_T^k \underline{\tau}_T, v)_T = -(\tau_T, \nabla v)_T + \sum_{F \in \mathcal{F}_T} (\tau_F \epsilon_{TF}, v)_F, \quad (11)$$

where $\epsilon_{TF} := \mathbf{n}_{TF} \cdot \mathbf{n}_F$ for all $T \in \mathcal{T}_h$ and all $F \in \mathcal{F}_T$.

Lemma 2 (Commuting property). *The following holds for all $\tau \in \Sigma^+(T)$:*

$$D_T^k(\underline{I}_T^k \tau) = \pi_T^k(\nabla \cdot \tau). \quad (12)$$

Proof. For all $v \in \mathbb{P}_d^k(T)$, we observe that

$$\begin{aligned} (\pi_T^k(\nabla \cdot \tau), v)_T &= (\nabla \cdot \tau, v)_T = -(\tau, \nabla v)_T + \sum_{F \in \mathcal{F}_T} (\tau \cdot \mathbf{n}_{TF}, v)_F \\ &= -((\underline{I}_T^k \tau)_T, \nabla v)_T + \sum_{F \in \mathcal{F}_T} ((\underline{I}_T^k \tau)_F \epsilon_{TF}, v)_F = (D_T^k(\underline{I}_T^k \tau), v)_T, \end{aligned}$$

where we have used integration by parts in T , the definition of \underline{I}_T^k as an element of $\underline{\Sigma}_T^k$, and that of D_T^k acting on an element of $\underline{\Sigma}_T^k$. \square

2.5 Consistent flux reconstruction

The consistent flux reconstruction operator $\mathfrak{C}_T^{k+1} : \underline{\Sigma}_T^k \rightarrow \mathbf{\Gamma}_T^{k+1} := \mathbf{K}_T \nabla \mathbb{P}_d^{k+1}(T)$ is such that, for all $\underline{\tau}_T \in \underline{\Sigma}_T^k$, $\mathfrak{C}_T^{k+1} \underline{\tau}_T = \mathbf{K}_T \nabla v$ where $v \in \mathbb{P}_d^{k+1}(T)$ solves the following Neumann problem: For all $w \in \mathbb{P}_d^{k+1}(T)$,

$$(\mathfrak{C}_T^{k+1} \underline{\tau}_T, \nabla w)_T = (\mathbf{K}_T \nabla v, \nabla w)_T = -(D_T^k \underline{\tau}_T, w)_T + \sum_{F \in \mathcal{F}_T} (\tau_F \epsilon_{TF}, w)_F, \quad (13)$$

with compatible right-hand side vanishing for constant w owing to (11) (while v is defined up to a constant which can be fixed prescribing its average value on T).

Lemma 3 (Polynomial preservation). *The following holds for all $\boldsymbol{\tau} \in \boldsymbol{\Gamma}_T^{k+1}$:*

$$\mathfrak{C}_T^{k+1}(\underline{I}_T^k \boldsymbol{\tau}) = \boldsymbol{\tau}. \quad (14)$$

Proof. Let $\boldsymbol{\tau} \in \boldsymbol{\Gamma}_T^{k+1}$. Owing to (13), we infer that, for all $w \in \mathbb{P}_d^{k+1}(T)$,

$$(\mathfrak{C}_T^{k+1}(\underline{I}_T^k \boldsymbol{\tau}), \nabla w)_T = -(D_T^k(\underline{I}_T^k \boldsymbol{\tau}), w)_T + \sum_{F \in \mathcal{F}_T} ((\underline{I}_T^k \boldsymbol{\tau})_F \epsilon_{TF}, w)_F.$$

The commuting property (12) implies that $D_T^k(\underline{I}_T^k \boldsymbol{\tau}) = \pi_T^k(\nabla \cdot \boldsymbol{\tau}) = \nabla \cdot \boldsymbol{\tau}$ since $\boldsymbol{\tau} \in \boldsymbol{\Gamma}_T^{k+1} \subset \mathbb{P}_d^k(T)$. For the same reason, $(\underline{I}_T^k \boldsymbol{\tau})_F = \pi_F^k(\boldsymbol{\tau} \cdot \mathbf{n}_F) = \boldsymbol{\tau} \cdot \mathbf{n}_F$. As a result,

$$(\mathfrak{C}_T^{k+1}(\underline{I}_T^k \boldsymbol{\tau}), \nabla w)_T = -(\nabla \cdot \boldsymbol{\tau}, w)_T + \sum_{F \in \mathcal{F}_T} (\boldsymbol{\tau} \cdot \mathbf{n}_{TF}, w)_F = (\boldsymbol{\tau}, \nabla w)_T,$$

which proves (14) since $(\mathfrak{C}_T^{k+1}(\underline{I}_T^k \boldsymbol{\tau}) - \boldsymbol{\tau}) \in \boldsymbol{\Gamma}_T^{k+1} = \mathbf{K}_T \nabla \mathbb{P}_d^{k+1}(T)$. \square

2.6 Discrete Hodge inner product

The discrete Hodge inner product $H_T : \underline{\Sigma}_T^k \times \underline{\Sigma}_T^k \rightarrow \mathbb{R}$ is such that, for all $\underline{\boldsymbol{\sigma}}_T, \underline{\boldsymbol{\tau}}_T \in \underline{\Sigma}_T^k$,

$$H_T(\underline{\boldsymbol{\sigma}}_T, \underline{\boldsymbol{\tau}}_T) := (\mathbf{K}_T^{-1} \mathfrak{C}_T^{k+1} \underline{\boldsymbol{\sigma}}_T, \mathfrak{C}_T^{k+1} \underline{\boldsymbol{\tau}}_T)_T + S_T(\underline{\boldsymbol{\sigma}}_T, \underline{\boldsymbol{\tau}}_T), \quad (15)$$

with stabilization bilinear form S_T such that, letting $\kappa_{TF} := \mathbf{n}_F \cdot \mathbf{K}_T \cdot \mathbf{n}_F$,

$$S_T(\underline{\boldsymbol{\sigma}}_T, \underline{\boldsymbol{\tau}}_T) := \sum_{F \in \mathcal{F}_T} h_F \kappa_{TF}^{-1} ((\mathfrak{C}_T^{k+1} \underline{\boldsymbol{\sigma}}_T) \cdot \mathbf{n}_F - \sigma_F, (\mathfrak{C}_T^{k+1} \underline{\boldsymbol{\tau}}_T) \cdot \mathbf{n}_F - \tau_F)_F. \quad (16)$$

Notice that the stabilization bilinear form is symmetric and positive semi-definite, so that introducing the semi-norm $|\underline{\boldsymbol{\tau}}_T|_{S,T} := S_T(\underline{\boldsymbol{\tau}}_T, \underline{\boldsymbol{\tau}}_T)^{1/2}$ on $\underline{\Sigma}_T^k$, we infer that

$$S_T(\underline{\boldsymbol{\sigma}}_T, \underline{\boldsymbol{\tau}}_T) \leq |\underline{\boldsymbol{\sigma}}_T|_{S,T} |\underline{\boldsymbol{\tau}}_T|_{S,T} \quad \forall \underline{\boldsymbol{\sigma}}_T, \underline{\boldsymbol{\tau}}_T \in \underline{\Sigma}_T^k. \quad (17)$$

Another important property of S_T is the following polynomial consistency: For all $\boldsymbol{\sigma} \in \boldsymbol{\Gamma}_T^{k+1}$,

$$S_T(\underline{I}_T^k(\boldsymbol{\sigma}), \underline{\boldsymbol{\tau}}_T) = 0 \quad \forall \underline{\boldsymbol{\tau}}_T \in \underline{\Sigma}_T^k. \quad (18)$$

This is a consequence of the fact that $\mathfrak{C}_T^{k+1}(\underline{I}_T^k \boldsymbol{\sigma}) = \boldsymbol{\sigma}$ owing to the polynomial consistency property (14) and that $(\underline{I}_T^k \boldsymbol{\sigma})_F = \boldsymbol{\sigma} \cdot \mathbf{n}_F$ since $\boldsymbol{\Gamma}_T^{k+1} \subset \mathbb{P}_d^k(T)$.

3 Discrete problem and main results

In this section, we formulate the discrete problem and state our main results; their proofs are postponed to Section 4. We also discuss static condensation, which is important in practice, and draw links with existing methods from the literature in the lowest-order case ($k = 0$).

3.1 Discrete problem

The global DoFs for the flux and the potential are

$$\underline{\Sigma}_h^k := \left\{ \bigtimes_{T \in \mathcal{T}_h} \mathbf{\Gamma}_T^k \right\} \times \left\{ \bigtimes_{F \in \mathcal{F}_h} \mathbb{P}_{d-1}^k(F) \right\}, \quad U_h^k := \bigtimes_{T \in \mathcal{T}_h} \mathbb{P}_d^k(T), \quad (19)$$

so that the face-based DoFs of the flux are patched. A generic collection of DoFs in $\underline{\Sigma}_h^k$ is denoted $\underline{\tau}_h = ((\tau_T)_{T \in \mathcal{T}_h}, (\tau_F)_{F \in \mathcal{F}_h})$, and for all $T \in \mathcal{T}_h$, we set $\underline{\tau}_T := (\tau_T, (\tau_F)_{F \in \mathcal{F}_T}) \in \underline{\Sigma}_T^k$. A generic element in U_h^k is denoted $v_h = (v_T)_{T \in \mathcal{T}_h}$.

The discrete problem consists in finding $(\underline{\sigma}_h, u_h) \in \underline{\Sigma}_h^k \times U_h^k$ such that, for all $(\underline{\tau}_h, v_h) \in \underline{\Sigma}_h^k \times U_h^k$, the following holds for all $T \in \mathcal{T}_h$:

$$H_T(\underline{\sigma}_T, \underline{\tau}_T) + (D_T^k \underline{\tau}_T, u_T)_T = 0, \quad (20a)$$

$$(D_T^k \underline{\sigma}_T, v_T)_T = -(f, v_T)_T. \quad (20b)$$

3.2 Stability and well-posedness

We introduce the following norms on $\underline{\Sigma}_T^k$:

$$\|\underline{\tau}_T\|_{H,T}^2 := H_T(\underline{\tau}_T, \underline{\tau}_T) = \|\mathbf{K}_T^{-1/2} \mathfrak{C}_T^{k+1} \underline{\tau}_T\|_T^2 + |\underline{\tau}_T|_{S,T}^2, \quad (21a)$$

$$\|\underline{\tau}_T\|_T^2 := \|\tau_T\|_T^2 + \sum_{F \in \mathcal{F}_T} h_F \|\tau_F\|_F^2. \quad (21b)$$

It is clear that $\|\cdot\|_T$ defines a norm on $\underline{\Sigma}_T^k$; that $\|\cdot\|_{H,T}$ also defines a norm follows from the following result.

Lemma 4 (Stability of H_T). *There is $\eta > 0$, uniform with respect to the mesh size and the diffusion tensor, such that the following holds:*

$$\eta \lambda_{\sharp,T}^{-1/2} \|\underline{\tau}_T\|_T \leq \|\underline{\tau}_T\|_{H,T} \leq \eta^{-1} \lambda_{\flat,T}^{-1/2} \|\underline{\tau}_T\|_T, \quad (22)$$

for all $T \in \mathcal{T}_h$ and all $\underline{\tau}_T \in \underline{\Sigma}_T^k$.

We introduce the global flux norm such that $\|\underline{\tau}_h\|_H^2 := \sum_{T \in \mathcal{T}_h} \|\underline{\tau}_T\|_{H,T}^2$ for all $\underline{\tau}_h \in \underline{\Sigma}_h^k$. The following result is a classical consequence of the above setting.

Lemma 5 (Well-posedness of (20)). *There exists a real number $\beta > 0$, uniform with respect to the mesh size and the diffusion tensor, such that, for all $v_h \in U_h^k$, the following holds:*

$$\lambda_{\flat}^{1/2} \beta \|v_h\| \leq \sup_{\underline{\tau}_h \in \underline{\Sigma}_h^k, \|\underline{\tau}_h\|_H=1} \left\{ \sum_{T \in \mathcal{T}_h} (D_T^k \underline{\tau}_T, v_T)_T \right\}. \quad (23)$$

Additionally, problem (20) is well-posed.

3.3 Error estimates

Assuming that $\sigma|_T \in \Sigma^+(T)$ for all $T \in \mathcal{T}_h$, we define the discrete objects $(\hat{\underline{\sigma}}_h, \hat{u}_h) \in \underline{\Sigma}_h^k \times U_h^k$ such that, for all $T \in \mathcal{T}_h$,

$$\hat{\underline{\sigma}}_T := \underline{I}_T^k(\sigma|_T), \quad \hat{u}_T := \pi_T^k(u|_T). \quad (24)$$

The definition of $\hat{\underline{\sigma}}_h$ is meaningful since $\sigma \cdot \mathbf{n}_F$ is single-valued for all $F \in \mathcal{F}_h$.

Theorem 6 (Flux-error estimate). *Let (σ, u) be the unique solution to (2) and let $(\underline{\sigma}_h, u_h)$ be the unique solution to (20). Assume that $u \in H^{k+2}(\mathcal{T}_h)$. Then, the following holds:*

$$\|\hat{\underline{\sigma}}_h - \underline{\sigma}_h\|_H \lesssim \left\{ \sum_{T \in \mathcal{T}_h} \rho_{\mathbf{K}, T} \lambda_{\sharp, T} h_T^{2(k+1)} |u|_{H^{k+2}(T)}^2 \right\}^{1/2}, \quad (25)$$

and

$$\left\{ \sum_{T \in \mathcal{T}_h} \|\mathbf{K}_T^{-1/2}(\mathfrak{C}_T^{k+1} \underline{\sigma}_T - \sigma)\|_T^2 \right\}^{1/2} \lesssim \left\{ \sum_{T \in \mathcal{T}_h} \rho_{\mathbf{K}, T} \lambda_{\sharp, T} h_T^{2(k+1)} |u|_{H^{k+2}(T)}^2 \right\}^{1/2}. \quad (26)$$

Defining the function $\hat{u}_h \in U_h^k$ such that $\hat{u}_h|_T = \hat{u}_T$ for all $T \in \mathcal{T}_h$, a potential L^2 -error estimate of order $(k+1)$ bounding $\|\hat{u}_h - u_h\|$ by the right-hand side of (25) follows from Lemma 5 and Theorem 6 (see Remark 10 below). An improved error estimate on the potential holds under the following elliptic regularity assumption: There is a real number $C_{\text{ell}} > 0$, only depending on Ω , such that, for all $g \in L^2(\Omega)$, the unique solution $z \in H_0^1(\Omega)$ of $-\nabla \cdot (\mathbf{K} \nabla z) = g$ satisfies $\|z\|_{H^2(\Omega)} \leq C_{\text{ell}} \lambda_{\flat}^{-1/2} \|g\|$.

Theorem 7 (Supercloseness of the potential). *Assume elliptic regularity and, for $k = 0$ that $f \in H^1(\mathcal{T}_h)$. Then, under the assumptions of Theorem 6, the following holds:*

$$\|\hat{u}_h - u_h\| \lesssim \rho_{\mathbf{K}} h \left\{ \sum_{T \in \mathcal{T}_h} \rho_{\mathbf{K}, T} \lambda_{\sharp, T} h_T^{2(k+1)} \|u\|_{H^{k+2}(T)}^2 \right\}^{1/2} + h^{k+2} \|f\|_{H^{k+\delta}(\mathcal{T}_h)}, \quad (27)$$

where $\rho_{\mathbf{K}} := \lambda_{\sharp}/\lambda_{\flat}$ while $\delta = 1$ for $k = 0$ and $\delta = 0$ for $k \geq 1$.

3.4 Static condensation

We briefly discuss two approaches for reducing substantially the size of the discrete problem (20) by means of static condensation, the second approach being more computationally effective.

In the first approach, we eliminate locally the cell-based flux DoFs and the potential DoFs up to one constant value per mesh cell. Let U_T^0 be spanned by constant potentials and let $U_T^{k,0}$ be spanned by polynomials of degree at most k having zero mean-value in T . Observe that $U_T^k = U_T^0 \oplus U_T^{k,0}$ and correspondingly write $u_T = (u_T^0, \tilde{u}_T)$ for the potential with $u_T^0 \in U_T^0$ and $\tilde{u}_T \in U_T^{k,0}$. Then, we infer from (20) that, for all $T \in \mathcal{T}_h$, $(\sigma_T, \tilde{u}_T) \in \mathbf{\Gamma}_T^k \times U_T^{k,0}$ can be eliminated locally by solving the following saddle-point problem:

$$\tilde{H}_T(\sigma_T, \tau_T) + (\tau_T, \nabla \tilde{u}_T)_T = g_1(\tau_T), \quad (28a)$$

$$(\sigma_T, \nabla \tilde{v}_T)_T = g_2(\tilde{v}_T), \quad (28b)$$

for all $(\boldsymbol{\tau}_T, \tilde{v}_T) \in (\mathbf{K}_T \nabla \mathbb{P}_d^k(T)) \times U_T^{k,0}$ where g_1, g_2 are suitable linear forms and $\tilde{H}_T(\boldsymbol{\sigma}_T, \boldsymbol{\tau}_T) := H_T((\boldsymbol{\sigma}_T, (0)_{F \in \mathcal{F}_T}), (\boldsymbol{\tau}_T, (0)_{F \in \mathcal{F}_T}))$. Owing to Lemma 4, we infer that $\tilde{H}_T(\boldsymbol{\tau}_T, \boldsymbol{\tau}_T)$ is uniformly equivalent to $\|\boldsymbol{\tau}_T\|_T^2$, so that (28) is well-posed. After static condensation, the global linear system is still of saddle-point form and involves the face-based flux DoFs and the mean-value of the potential in each mesh cell. This problem has the same size and structure as that derived in the MHM method [2, 39].

The second approach is closely inspired by the hybridization technique for mixed finite elements introduced in [3]. The key idea consists in enforcing the single-valuedness of interface flux unknowns on every $F \in \mathcal{F}_h^i$ by means of Lagrange multipliers in $\mathbb{P}_{d-1}^k(F)$, thereafter recovering a primal problem once (cell- and face-) flux unknowns have been locally eliminated. Note that the Lagrange multipliers can be interpreted as potential traces. Let $T \in \mathcal{T}_h$ and consider a local collection of potential DoFs and Lagrange multipliers

$$\underline{v}_T = (v_T, (v_F)_{F \in \mathcal{F}_T}) \in U_T^k \times \left\{ \bigtimes_{F \in \mathcal{F}_T} \Lambda_F^k \right\} =: \underline{W}_T^k,$$

with $\Lambda_F^k := \mathbb{P}_{d-1}^k(F)$ if $F \in \mathcal{F}_T \cap \mathcal{F}_h^i$ while $\Lambda_F^k := \{0\}$ if $F \in \mathcal{F}_T \cap \mathcal{F}_h^b$. To eliminate the flux unknowns in $\underline{\Sigma}_T^k$, we introduce the local operator $\underline{\varsigma}_T : \underline{W}_T^k \rightarrow \underline{\Sigma}_T^k$ such that, for all $\underline{v}_T \in \underline{W}_T^k$, $\underline{\varsigma}_T(\underline{v}_T) \in \underline{\Sigma}_T^k$ solves the following local problem:

$$H_T(\underline{\varsigma}_T(\underline{v}_T), \boldsymbol{\tau}_T) = -(D_T^k \boldsymbol{\tau}_T, v_T)_T + \sum_{F \in \mathcal{F}_T} (\tau_F \epsilon_{TF}, v_F)_F \quad \forall \boldsymbol{\tau}_T \in \underline{\Sigma}_T^k. \quad (29)$$

The well-posedness of (29) classically follows from (22). Define now the global space of potential unknowns and Lagrange multipliers as

$$\underline{W}_h^k := U_h^k \times \left\{ \bigtimes_{F \in \mathcal{F}_h} \Lambda_F^k \right\}.$$

Denoting by $(\boldsymbol{\sigma}_h, \underline{u}_h) \in \underline{\Sigma}_h^k \times \underline{W}_h^k$ the unique solution to the problem obtained from (20) by enforcing the single-valuedness of face unknowns for the flux via Lagrange multipliers, one can easily show that $\boldsymbol{\sigma}_T = \underline{\varsigma}_T(\underline{u}_T)$ for all $T \in \mathcal{T}_h$. Additionally, \underline{u}_h can be obtained independently from $\boldsymbol{\sigma}_h$ by solving the following primal problem where it appears as the sole unknown: Find $\underline{u}_h \in \underline{W}_h^k$ such that, for all $\underline{v}_h \in \underline{W}_h^k$, the following holds for all $T \in \mathcal{T}_h$:

$$H_T(\underline{\varsigma}_T(\underline{u}_T), \underline{\varsigma}_T(\underline{v}_T)) = (f, v_T)_T. \quad (30)$$

The well-posedness of (30) stems from the stability (22) of H_T and the injectivity of $\underline{\varsigma}_T$. Conversely, if $\underline{u}_h \in \underline{W}_h^k$ solves (30), then setting $\boldsymbol{\sigma}_T = \underline{\varsigma}_T(\underline{u}_T) \in \underline{\Sigma}_T^k$ for all $T \in \mathcal{T}_h$, observing that σ_F is single-valued at all mesh interfaces so that $\boldsymbol{\sigma}_h \in \underline{\Sigma}_h^k$, and letting $u_h \in U_h^k$ collect the cell DoFs of \underline{u}_h , one can prove that the pair $(\boldsymbol{\sigma}_h, u_h)$ solves (20). As a result, one can solve the coercive primal problem (30) in place of the saddle-point problem (20). Additionally, the size of the global system in (30) can be further reduced by performing static condensation to express the (cell) potential unknowns in terms of the Lagrange multipliers. Both flux unknowns and (cell) potential unknowns can then be recovered by local post-processing. As a closing remark, we observe that the primal problem (30) has the same structure as the Hybrid High-Order method of [31].

3.5 Lowest-order case ($k = 0$)

Since lowest-order mixed methods have been extensively explored in the literature, we devote this section to a brief discussion of the links with the present method in the case where $k = 0$ recalling that only the face-based flux DoFs are relevant, i.e., $\underline{\tau}_T = (\tau_F)_{F \in \mathcal{F}_T}$ for all $\underline{\tau}_T \in \underline{\Sigma}_T^0$. We first observe that (13) leads to the following explicit expression:

$$\mathfrak{C}_T^1 \underline{\tau}_T = \frac{1}{|T|_d} \sum_{F \in \mathcal{F}_T} |F|_{d-1} (\bar{\mathbf{x}}_F - \bar{\mathbf{x}}_T) \epsilon_{TF} \tau_F, \quad (31)$$

where $\bar{\mathbf{x}}_F$ and $\bar{\mathbf{x}}_T$ denote the barycenter of F and T , respectively. This lowest-order explicit flux reconstruction is well-known in the context of Mixed Finite Volumes; see [34, eq. (9)].

To draw further links, let us observe that the stabilization bilinear form can be taken in the non-diagonal form

$$S_T(\underline{\sigma}_T, \underline{\tau}_T) = \sum_{F, F' \in \mathcal{F}_T} ((\mathfrak{C}_T^1 \underline{\sigma}_T) \cdot \mathbf{n}_F - \sigma_F) \mathbb{M}_{T, FF'} ((\mathfrak{C}_T^1 \underline{\tau}_T) \cdot \mathbf{n}_{F'} - \tau_{F'}), \quad (32)$$

where \mathbb{M}_T is a symmetric positive definite matrix of order $\#(\mathcal{F}_T)$. The choice (16) corresponds to a diagonal matrix with diagonal entries set to $|F|_{d-1} h_F \kappa_{TF}^{-1}$ for all $F \in \mathcal{F}_T$. Another approach to design the matrix \mathbb{M}_T is to use reconstruction functions $\{\varphi_{TF}\}_{F \in \mathcal{F}_T}$. Denoting $\bar{\varphi}$ the mean-value of a generic function φ in T , the reconstruction functions must satisfy

$$\bar{\varphi}_{TF} = \frac{|F|_{d-1}}{|T|_d} (\bar{\mathbf{x}}_F - \bar{\mathbf{x}}_T), \quad \sum_{F \in \mathcal{F}_T} \varphi_{TF}(\mathbf{x}) \otimes \mathbf{n}_{TF} = \mathbf{I}_d, \quad (33)$$

where \mathbf{I}_d is the identity matrix in $\mathbb{R}^{d \times d}$. The first property in (33) implies that the reconstruction operator $\mathfrak{R}_T(\underline{\tau}_T)(\mathbf{x}) := \sum_{F \in \mathcal{F}_T} \tau_F \epsilon_{TF} \varphi_{TF}(\mathbf{x})$ is such that

$$\overline{\mathfrak{R}_T(\underline{\tau}_T)} = \mathfrak{C}_T^1(\underline{\tau}_T). \quad (34)$$

The second property implies that

$$\sum_{F \in \mathcal{F}_T} (\mathfrak{C}_T^1(\underline{\tau}_T) \cdot \mathbf{n}_F) \epsilon_{TF} \varphi_{TF}(\mathbf{x}) = \mathfrak{C}_T^1(\underline{\tau}_T) = \overline{\mathfrak{C}_T^1(\underline{\tau}_T)} = \sum_{F \in \mathcal{F}_T} (\mathfrak{C}_T^1(\underline{\tau}_T) \cdot \mathbf{n}_F) \epsilon_{TF} \bar{\varphi}_{TF}. \quad (35)$$

As a result, defining the discrete Hodge inner product such that

$$H_T(\underline{\sigma}_T, \underline{\tau}_T) = \int_T \mathfrak{R}_T(\underline{\sigma}_T)(\mathbf{x}) \cdot \mathbf{K}_T^{-1} \cdot \mathfrak{R}_T(\underline{\tau}_T)(\mathbf{x}) d\mathbf{x}, \quad (36)$$

we infer that we recover definition (15) whenever the stabilization bilinear form results from (32) with non-diagonal matrix \mathbb{M}_T having entries given by

$$\mathbb{M}_{T, FF'} = \int_T (\varphi_{TF}(\mathbf{x}) - \bar{\varphi}_{TF}) \cdot \mathbf{K}_T^{-1} \cdot (\varphi_{TF'}(\mathbf{x}) - \bar{\varphi}_{TF'}) d\mathbf{x}. \quad (37)$$

Whenever \mathbf{K}_T is isotropic, this matrix is uniformly equivalent to the diagonal matrix associated with (16). Defining the discrete Hodge inner product as in (36) provides a link with CDO schemes in cell-based form [12]. Examples of reconstruction functions are the Raviart–Thomas basis functions on simplices and the piecewise constant (on a simplicial submesh) functions on polyhedral cells from the Discrete Geometric Approach [26].

4 Proofs

This section collects the proofs of our results.

4.1 Preliminary results

Lemma 8 (Stability of D_T^k , \mathfrak{C}_T^{k+1} , and S_T). *The following holds:*

$$h_T \|D_T^k \underline{\tau}_T\|_T + \lambda_{b,T}^{1/2} \|\mathbf{K}_T^{-1/2} \mathfrak{C}_T^{k+1} \underline{\tau}_T\|_T + \lambda_{b,T}^{1/2} |\underline{\tau}_T|_{S,T} \lesssim \|\underline{\tau}_T\|_T, \quad (38)$$

for all $T \in \mathcal{T}_h$ and all $\underline{\tau}_T \in \underline{\Sigma}_T^k$.

Proof. That $\|D_T^k \underline{\tau}_T\|_T \lesssim h_T^{-1} \|\underline{\tau}_T\|_T$ follows from the definition (11) of $D_T^k \underline{\tau}_T$ followed by an inverse inequality on ∇v , and a discrete trace inequality on $v|_F$ for all $F \in \mathcal{F}_T$. To bound $\|\mathbf{K}_T^{-1/2} \mathfrak{C}_T^{k+1} \underline{\tau}_T\|_T$, we write $\mathfrak{C}_T^{k+1} \underline{\tau}_T = \mathbf{K}_T \nabla v$ for some $v \in \mathbb{P}_d^{k+1}(T)$, and use the definition (13) of $\mathfrak{C}_T^{k+1} \underline{\tau}_T$ to infer that

$$\|\mathbf{K}_T^{-1/2} \mathfrak{C}_T^{k+1} \underline{\tau}_T\|_T^2 = (\mathfrak{C}_T^{k+1} \underline{\tau}_T, \nabla v)_T = -(D_T^k \underline{\tau}_T, \pi_T^k v)_T + \sum_{F \in \mathcal{F}_T} (\tau_F \epsilon_{TF}, \pi_F^k v)_F.$$

Since $(D_T^k \underline{\tau}_T, 1)_T = \sum_{F \in \mathcal{F}_T} (\tau_F \epsilon_{TF}, 1)_F$ owing to (11), we can write

$$\|\mathbf{K}_T^{-1/2} \mathfrak{C}_T^{k+1} \underline{\tau}_T\|_T^2 = -(D_T^k \underline{\tau}_T, \pi_T^k v - \pi_T^0 v)_T + \sum_{F \in \mathcal{F}_T} (\tau_F \epsilon_{TF}, \pi_F^k v - \pi_T^0 v)_F.$$

Since $\|\pi_T^k v - \pi_T^0 v\|_T \lesssim h_T \|\nabla v\|_T$, $\|\pi_F^k v - \pi_T^0 v\|_F \lesssim h_F^{1/2} \|\nabla v\|_T$, and $\|\nabla v\|_T \leq \lambda_{b,T}^{-1/2} \|\mathbf{K}_T^{1/2} \nabla v\|_T$, we infer from the above bound on $\|D_T^k \underline{\tau}_T\|_T$ that $\|\mathbf{K}_T^{-1/2} \mathfrak{C}_T^{k+1} \underline{\tau}_T\|_T \lesssim \lambda_{b,T}^{-1/2} \|\underline{\tau}_T\|_T$. Finally, to bound $|\underline{\tau}_T|_{S,T}$, we first observe that

$$|(\mathfrak{C}_T^{k+1} \underline{\tau}_T) \cdot \mathbf{n}_F| = |\mathbf{n}_F \cdot \mathbf{K}_T \cdot (\mathbf{K}_T^{-1} \mathfrak{C}_T^{k+1} \underline{\tau}_T)| \leq \kappa_{TF}^{1/2} |(\mathbf{K}_T^{-1} \mathfrak{C}_T^{k+1} \underline{\tau}_T) \cdot (\mathfrak{C}_T^{k+1} \underline{\tau}_T)|^{1/2},$$

since for two vectors $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$, $|\mathbf{x} \cdot \mathbf{K}_T \cdot \mathbf{y}| \leq |\mathbf{x} \cdot \mathbf{K}_T \cdot \mathbf{x}|^{1/2} |\mathbf{y} \cdot \mathbf{K}_T \cdot \mathbf{y}|^{1/2}$. As a result,

$$\|(\mathfrak{C}_T^{k+1} \underline{\tau}_T) \cdot \mathbf{n}_F\|_F \leq \kappa_{TF}^{1/2} \|\mathbf{K}_T^{-1/2} \mathfrak{C}_T^{k+1} \underline{\tau}_T\|_F. \quad (39)$$

Hence, using a triangle inequality, a discrete trace inequality to bound $\|\mathbf{K}_T^{-1/2} \mathfrak{C}_T^{k+1} \underline{\tau}_T\|_F$, and the above bound on $\|\mathbf{K}_T^{-1/2} \mathfrak{C}_T^{k+1} \underline{\tau}_T\|_T$, we infer that

$$\|(\mathfrak{C}_T^{k+1} \underline{\tau}_T) \cdot \mathbf{n}_F - \tau_F\|_F \lesssim h_T^{-1/2} \kappa_{TF}^{1/2} \lambda_{b,T}^{-1/2} \|\underline{\tau}_T\|_T + \|\tau_F\|_F,$$

whence the bound $|\underline{\tau}_T|_{S,T} \lesssim \lambda_{b,T}^{-1/2} \|\underline{\tau}_T\|_T$ follows from mesh regularity and the fact that $\kappa_{TF}^{-1/2} \leq \lambda_{b,T}^{-1/2}$ for all $F \in \mathcal{F}_T$. \square

Lemma 9 (Approximation properties of \mathfrak{C}_T^{k+1} and S_T). *For all $T \in \mathcal{T}_h$ and all $v \in H^{k+2}(T)$, letting $\tau := \mathbf{K}_T \nabla v$, the following holds:*

$$\|\mathbf{K}_T^{-1/2} (\mathfrak{C}_T^{k+1} (\underline{I}_T^k \tau) - \tau)\|_T + S_T (\underline{I}_T^k \tau, \underline{I}_T^k \tau)^{1/2} \lesssim \rho_{\mathbf{K},T}^{1/2} \lambda_{\sharp,T}^{1/2} h_T^{k+1} |v|_{H^{k+2}(T)}. \quad (40)$$

Proof. (1) *Bound on $\|\mathbf{K}_T^{-1/2}(\mathfrak{C}_T^{k+1}(\underline{I}_T^k \boldsymbol{\tau}) - \boldsymbol{\tau})\|_T$.* Let $\check{v}_T \in \mathbb{P}_d^{k+1}(T)$ solve the following well-posed Neumann problem:

$$(\mathbf{K}_T \nabla \check{v}_T, \nabla w)_T = (\mathbf{K}_T \nabla v, \nabla w)_T \quad \forall w \in \mathbb{P}_d^{k+1}(T), \quad (41)$$

with $\int_T \check{v}_T = \int_T v$. Using the triangle inequality, we infer that

$$\|\mathbf{K}_T^{-1/2}(\mathfrak{C}_T^{k+1}(\underline{I}_T^k \boldsymbol{\tau}) - \boldsymbol{\tau})\|_T \leq \|\mathbf{K}_T^{1/2} \nabla(\check{v}_T - v)\|_T + \|\mathbf{K}_T^{-1/2} \mathfrak{C}_T^{k+1}(\underline{I}_T^k \boldsymbol{\tau}) - \mathbf{K}_T^{1/2} \nabla \check{v}_T\|_T. \quad (42)$$

By definition, \check{v}_T is the element of $\mathbb{P}_d^{k+1}(T)$ which minimizes the distance to v in the $\|\mathbf{K}_T^{1/2} \nabla \cdot\|_T$ -norm, hence we can estimate the first term in (42) as follows:

$$\|\mathbf{K}_T^{1/2} \nabla(\check{v}_T - v)\|_T \leq \|\mathbf{K}_T^{1/2} \nabla(\pi_T^{k+1} v - v)\|_T \lesssim \lambda_{\sharp, T}^{1/2} h_T^{k+1} |v|_{H^{k+2}(T)}, \quad (43)$$

where we have concluded using the approximation properties (7) of π_T^{k+1} . Let us estimate the second term in the right-hand side of (42). Using the definitions (13) of \mathfrak{C}_T^{k+1} and (41) of \check{v}_T as well as that of the reduction operator \underline{I}_T^k and the commuting property (12) of D_T^k , we infer, for all $w \in \mathbb{P}_d^{k+1,0}(T)$ (the space of polynomials of degree $\leq k+1$ with zero average on T) that

$$(\mathfrak{C}_T^{k+1}(\underline{I}_T^k \boldsymbol{\tau}) - \mathbf{K}_T \nabla \check{v}_T, \nabla w)_T = (\nabla \cdot \boldsymbol{\tau} - \pi_T^k(\nabla \cdot \boldsymbol{\tau}), w)_T + \sum_{F \in \mathcal{F}_T} (\pi_F^k(\boldsymbol{\tau} \cdot \mathbf{n}_F) - \boldsymbol{\tau} \cdot \mathbf{n}_F, w|_F)_F.$$

Denote by \mathfrak{T}_1 and \mathfrak{T}_2 the addends in the right-hand side. For the first term, using the Cauchy–Schwarz inequality followed by the approximation properties (7) of π_T^{k+1} and the Poincaré inequality (8), we infer that $|\mathfrak{T}_1| \lesssim h_T^k |\nabla \cdot \boldsymbol{\tau}|_{H^k(T)} h_T \|\nabla w\|_T \lesssim \lambda_{\sharp, T} h_T^{k+1} |v|_{H^{k+2}(T)} \|\nabla w\|_T$. For the second term, we use the Cauchy–Schwarz inequality, the approximation properties (7) of the L^2 -orthogonal projector, the discrete trace inequality (5) and the Poincaré inequality (8) to infer $|\mathfrak{T}_2| \lesssim \lambda_{\sharp, T} h_T^{k+1} |v|_{H^{k+2}(T)} \|\nabla w\|_T$. Then, collecting the above bounds and since $\|\mathbf{K}_T^{-1/2} \boldsymbol{\tau}\|_T = \sup_{w \in \mathbb{P}_d^{k+1,0}(T)} \frac{(\boldsymbol{\tau}, \nabla w)_T}{\|\mathbf{K}_T^{1/2} \nabla w\|_T}$ for all $\boldsymbol{\tau} \in \boldsymbol{\Gamma}_T^{k+1}$, we infer that

$$\|\mathbf{K}_T^{-1/2} \mathfrak{C}_T^{k+1}(\underline{I}_T^k \boldsymbol{\tau}) - \mathbf{K}_T^{1/2} \nabla \check{v}_T\|_T \lesssim \rho_{\mathbf{K}, T}^{1/2} \lambda_{\sharp, T}^{1/2} h_T^{k+1} |v|_{H^{k+2}(T)}. \quad (44)$$

Using this bound and (43) in (42) together with $\rho_{\mathbf{K}, T} \geq 1$, the desired bound follows.

(2) *Bound on $S_T(\underline{I}_T^k \boldsymbol{\tau}, \underline{I}_T^k \boldsymbol{\tau})^{1/2}$.* We observe that, for all $F \in \mathcal{F}_T$,

$$\begin{aligned} \|\mathfrak{C}_T^{k+1}(\underline{I}_T^k \boldsymbol{\tau}) \cdot \mathbf{n}_F - \pi_F^k(\boldsymbol{\tau} \cdot \mathbf{n}_F)\|_F &= \|\pi_F^k((\mathfrak{C}_T^{k+1}(\underline{I}_T^k \boldsymbol{\tau}) - \boldsymbol{\tau}) \cdot \mathbf{n}_F)\|_F \\ &\leq \|(\mathfrak{C}_T^{k+1}(\underline{I}_T^k \boldsymbol{\tau}) - \boldsymbol{\tau}) \cdot \mathbf{n}_F\|_F \\ &\leq \kappa_{TF}^{1/2} \|\mathbf{K}_T^{-1/2}(\mathfrak{C}_T^{k+1}(\underline{I}_T^k \boldsymbol{\tau}) - \boldsymbol{\tau})\|_F, \end{aligned}$$

where we have used that $\mathfrak{C}_T^{k+1}(\underline{I}_T^k \boldsymbol{\tau}) \cdot \mathbf{n}_F \in \mathbb{P}_{d-1}^k(F)$, the fact that π_F^k is a projector, and a reasoning similar to the proof of (39). Adding and subtracting $\mathbf{K}_T \nabla \check{v}_T$ yields

$$\kappa_{TF}^{-1/2} \|\mathfrak{C}_T^{k+1}(\underline{I}_T^k \boldsymbol{\tau}) \cdot \mathbf{n}_F - \pi_F^k(\boldsymbol{\tau} \cdot \mathbf{n}_F)\|_F \leq \|\mathbf{K}_T^{-1/2} \mathfrak{C}_T^{k+1}(\underline{I}_T^k \boldsymbol{\tau}) - \mathbf{K}_T^{1/2} \nabla \check{v}_T\|_F + \|\mathbf{K}_T^{1/2} \nabla(\check{v}_T - v)\|_F.$$

We bound the first term in the right-hand side using a discrete trace inequality and (44), while we bound the second term by $\lambda_{\sharp, T}^{1/2} h_T^{k+1/2} |v|_{H^{k+2}(T)}$ using a continuous trace inequality followed by the Poincaré inequality (8) and (43). Since $\rho_{\mathbf{K}, T} \geq 1$, we infer that

$$\kappa_{TF}^{-1/2} \|\mathfrak{C}_T^{k+1}(\underline{I}_T^k \boldsymbol{\tau}) \cdot \mathbf{n}_F - \pi_F^k(\boldsymbol{\tau} \cdot \mathbf{n}_F)\|_F \lesssim \rho_{\mathbf{K}, T}^{1/2} \lambda_{\sharp, T}^{1/2} h_T^{k+1/2} |v|_{H^{k+2}(T)}.$$

Finally, the desired bound follows from the definition (16) of S_T and mesh regularity. \square

4.2 Proof of Lemma 4

Let $\underline{\tau}_T \in \underline{\Sigma}_T^k$ with $\underline{\tau}_T = (\tau_T, (\tau_F)_{F \in \mathcal{F}_T})$.

(1) *Lower bound on $\|\underline{\tau}_T\|_{H,T}$.* We write $\tau_T = \mathbf{K}_T \nabla v$ for some $v \in \mathbb{P}_d^k(T)$. Using (11) with test function v followed by (13) with test function v (this is possible since $\mathbb{P}_d^k(T) \subset \mathbb{P}_d^{k+1}(T)$), we infer that

$$\begin{aligned} (\mathbf{K}_T^{-1} \tau_T, \tau_T)_T &= (\tau_T, \nabla v)_T = -(D_T^k \underline{\tau}_T, v)_T + \sum_{F \in \mathcal{F}_T} (\tau_F \epsilon_{TF}, v)_F \\ &= (\mathfrak{C}_T^{k+1} \underline{\tau}_T, \nabla v)_T = (\mathbf{K}_T^{-1} \mathfrak{C}_T^{k+1} \underline{\tau}_T, \tau_T)_T. \end{aligned}$$

This identity readily implies that $\|\mathbf{K}_T^{-1/2} \tau_T\|_T \leq \|\mathbf{K}_T^{-1/2} \mathfrak{C}_T^{k+1} \underline{\tau}_T\|_T$, whence we infer that

$$\|\mathbf{K}_T^{-1/2} \mathfrak{C}_T^{k+1} \underline{\tau}_T\|_T \geq \lambda_{\sharp, T}^{-1/2} \|\tau_T\|_T. \quad (45)$$

Moreover, owing to the triangle inequality, we infer that

$$h_F^{1/2} \|\tau_F\|_F \leq h_F^{1/2} \|(\mathfrak{C}_T^{k+1} \underline{\tau}_T) \cdot \mathbf{n}_F - \tau_F\|_F + h_F^{1/2} \|(\mathfrak{C}_T^{k+1} \underline{\tau}_T) \cdot \mathbf{n}_F\|_F.$$

Using (39) followed by a discrete trace inequality to bound $\|\mathbf{K}_T^{-1/2} \mathfrak{C}_T^{k+1} \underline{\tau}_T\|_F$ yields

$$h_F^{1/2} \|\tau_F\|_F \lesssim h_F^{1/2} \|(\mathfrak{C}_T^{k+1} \underline{\tau}_T) \cdot \mathbf{n}_F - \tau_F\|_F + \kappa_{TF}^{1/2} \|\mathbf{K}_T^{-1/2} \mathfrak{C}_T^{k+1} \underline{\tau}_T\|_T.$$

Recalling the definition of $|\cdot|_{S,T}$, squaring and summing over $F \in \mathcal{F}_T$, we obtain

$$\sum_{F \in \mathcal{F}_T} h_F \|\tau_F\|_F^2 \lesssim \lambda_{\sharp, T}^2 |\underline{\tau}_T|_{S,T}^2 + \lambda_{\sharp, T} \|\mathbf{K}_T^{-1/2} \mathfrak{C}_T^{k+1} \underline{\tau}_T\|_T^2,$$

since $\kappa_{TF} \leq \lambda_{\sharp, T}$ for all $F \in \mathcal{F}_T$. Combining this bound with (45), we infer the desired lower bound on $\|\underline{\tau}_T\|_{H,T}$.

(2) *Upper bound on $\|\underline{\tau}_T\|_{H,T}$.* This bound is a straightforward consequence of (38).

4.3 Proof of Theorem 6

We start by observing that the following holds with local consistency error $\mathcal{E}_T(\underline{\tau}_T) := H_T(\hat{\underline{\sigma}}_T, \underline{\tau}_T) + (D_T^k \underline{\tau}_T, \hat{u}_T)_T$ for all $T \in \mathcal{T}_h$:

$$\|\hat{\underline{\sigma}}_h - \underline{\sigma}_h\|_H \leq \sup_{\underline{\tau}_h \in \underline{\Sigma}_h^k, \|\underline{\tau}_h\|_H = 1} \left\{ \sum_{T \in \mathcal{T}_h} \mathcal{E}_T(\underline{\tau}_T) \right\}. \quad (46)$$

Indeed, let $\underline{\tau}_h \in \underline{\Sigma}_h^k$ be such that $\|\underline{\tau}_h\|_H = 1$. Owing to (20a), we infer that

$$H_T(\hat{\underline{\sigma}}_T - \underline{\sigma}_T, \underline{\tau}_T) + (D_T^k \underline{\tau}_T, \hat{u}_T - u_T)_T = \mathcal{E}_T(\underline{\tau}_T).$$

Letting $\underline{\tau}_h = \frac{1}{\|\hat{\underline{\sigma}}_h - \underline{\sigma}_h\|_H} (\hat{\underline{\sigma}}_h - \underline{\sigma}_h)$ and since $D_T^k(\hat{\underline{\sigma}}_T - \underline{\sigma}_T) = 0$, for all $T \in \mathcal{T}_h$, owing to the commuting property (12) and the discrete equation (20b), the bound (46) follows.

To prove (25), we estimate $\mathcal{E}_T(\underline{\tau}_T)$ for all $T \in \mathcal{T}_h$ and all $\underline{\tau}_h \in \underline{\Sigma}_h^k$ such that $\|\underline{\tau}_h\|_H = 1$. We introduce the discrete functions $\check{u}_T := \pi_T^{k+1}(u|_T)$ and $\check{\underline{\sigma}}_T := \underline{I}_T^k(\mathbf{K}_T \nabla \check{u}_T)$ and decompose the local consistency error as follows:

$$\begin{aligned} \mathcal{E}_T(\underline{\tau}_T) &= H_T(\hat{\underline{\sigma}}_T - \check{\underline{\sigma}}_T, \underline{\tau}_T) + (D_T^k \underline{\tau}_T, \hat{u}_T - \check{u}_T)_T + \left\{ H_T(\check{\underline{\sigma}}_T, \underline{\tau}_T) + (D_T^k \underline{\tau}_T, \check{u}_T)_T \right\} \\ &:= \mathfrak{T}_{1,T} + \mathfrak{T}_{2,T} + \mathfrak{T}_{3,T}. \end{aligned}$$

(1) *Bound on $\mathfrak{I}_{1,T}$.* Since the discrete Hodge inner product is a symmetric and positive definite bilinear form, we infer that

$$\mathfrak{I}_{1,T} \leq \|\hat{\underline{\sigma}}_T - \check{\underline{\sigma}}_T\|_{H,T} \|\underline{\tau}_T\|_{H,T}.$$

Using the upper bound in (22) yields

$$\lambda_{b,T}^{1/2} \|\hat{\underline{\sigma}}_T - \check{\underline{\sigma}}_T\|_{H,T} \lesssim \|\hat{\underline{\sigma}}_T - \check{\underline{\sigma}}_T\|_T + \sum_{F \in \mathcal{F}_T} h_F^{1/2} \|\hat{\underline{\sigma}}_F - \check{\underline{\sigma}}_F\|_F,$$

where $\hat{\underline{\sigma}}_T$ and $\check{\underline{\sigma}}_T$ are the components in $\mathbf{\Gamma}_T^k$ of $\hat{\underline{\sigma}}_T$ and $\check{\underline{\sigma}}_T$, respectively, and $\hat{\underline{\sigma}}_F$ and $\check{\underline{\sigma}}_F$ the components in $\mathbb{P}_{d-1}^k(F)$. Recalling the definition of $\underline{\tau}_T^k$, see (10), we infer that

$$\hat{\underline{\sigma}}_T = \varpi_T^k(\mathbf{K}_T \nabla u), \quad \check{\underline{\sigma}}_T = \varpi_T^k(\mathbf{K}_T \nabla \check{u}_T),$$

where ϖ_T^k denotes the $(\mathbf{K}_T^{-1} \cdot, \cdot)_T$ -orthogonal projector onto $\mathbf{\Gamma}_T^k$. Since

$$\|\varpi_T^k \tau\|_T \leq \lambda_{\sharp,T}^{1/2} \|\mathbf{K}_T^{-1/2} \varpi_T^k \tau\|_T \leq \lambda_{\sharp,T}^{1/2} \|\mathbf{K}_T^{-1/2} \tau\|_T \leq \lambda_{\sharp,T} \|\mathbf{K}_T^{-1} \tau\|_T,$$

for all $\tau \in \mathbf{L}^2(T)$, using the approximation property (7) of π_T^{k+1} , we infer that

$$\|\hat{\underline{\sigma}}_T - \check{\underline{\sigma}}_T\|_T = \|\varpi_T^k(\mathbf{K}_T \nabla(u - \check{u}_T))\|_T \lesssim \lambda_{\sharp,T} h_T^{k+1} |u|_{H^{k+2}(T)}.$$

Moreover, $\hat{\underline{\sigma}}_F = \pi_F^k(\mathbf{K}_T \nabla u \cdot \mathbf{n}_F)$ and $\check{\underline{\sigma}}_F = \pi_F^k(\mathbf{K}_T \nabla \check{u}_T \cdot \mathbf{n}_F)$, so that

$$\|\hat{\underline{\sigma}}_F - \check{\underline{\sigma}}_F\|_F \leq \|\mathbf{K}_T \nabla(u - \check{u}_T) \cdot \mathbf{n}_F\|_F \lesssim \lambda_{\sharp,T} h_T^{k+1/2} |u|_{H^{k+2}(T)}.$$

Collecting the above bounds yields

$$\mathfrak{I}_{1,T} \lesssim \rho_{\mathbf{K},T}^{1/2} \lambda_{\sharp,T}^{1/2} h_T^{k+1} |u|_{H^{k+2}(T)} \|\underline{\tau}_T\|_{H,T}.$$

(2) *Bound on $\mathfrak{I}_{2,T}$.* We observe that

$$\begin{aligned} \mathfrak{I}_{2,T} &= (\hat{u}_T - \check{u}_T, D_T^k \underline{\tau}_T)_T = (\pi_T^k(u - \pi_T^{k+1} u), D_T^k \underline{\tau}_T)_T \leq \|u - \pi_T^{k+1} u\|_T \|D_T^k \underline{\tau}_T\|_T \\ &\lesssim \lambda_{\sharp,T}^{1/2} h_T^{k+1} |u|_{H^{k+2}(T)} \|\underline{\tau}_T\|_{H,T}, \end{aligned}$$

where we have used that $h_T \|D_T^k \underline{\tau}_T\|_T \lesssim \|\underline{\tau}_T\|_T \lesssim \lambda_{\sharp,T}^{1/2} \|\underline{\tau}_T\|_{H,T}$ owing to (38) and the lower bound in (22).

(3) *Reformulation of $\mathfrak{I}_{3,T}$.* Since $\mathbf{K}_T \nabla \check{u}_T \in \mathbf{\Gamma}_T^{k+1}$, we infer from (18) that

$$S_T(\check{\underline{\sigma}}_T, \underline{\tau}_T) = 0.$$

Moreover, (14) implies that $\mathfrak{C}_T^{k+1} \check{\underline{\sigma}}_T = \mathbf{K}_T \nabla \check{u}_T$. Hence,

$$\mathfrak{I}_{3,T} = (\nabla \check{u}_T, \mathfrak{C}_T^{k+1} \underline{\tau}_T)_T + (D_T^k \underline{\tau}_T, \check{u}_T)_T = \sum_{F \in \mathcal{F}_T} (\tau_F \epsilon_{TF}, \check{u}_T)_F,$$

where we have used (13) for the definition of $\mathfrak{C}_T^{k+1} \underline{\tau}_T$.

(4) *Conclusion.* We need to bound $\sum_{T \in \mathcal{T}_h} \{\mathfrak{I}_{1,T} + \mathfrak{I}_{2,T} + \mathfrak{I}_{3,T}\}$. Using Steps (1) and (2) and a discrete Cauchy–Schwarz inequality shows that $\sum_{T \in \mathcal{T}_h} \{\mathfrak{I}_{1,T} + \mathfrak{I}_{2,T}\}$ is bounded by the

right-hand side of (25). Furthermore, since the exact potential u is single-valued at interfaces and vanishes at boundary faces, we infer that

$$\sum_{T \in \mathcal{T}_h} \mathfrak{T}_{3,T} = \sum_{T \in \mathcal{T}_h} \sum_{F \in \mathcal{F}_T} (\tilde{u}_T - u, \tau_F \epsilon_{TF})_F,$$

so that $\sum_{T \in \mathcal{T}_h} \mathfrak{T}_{3,T}$ is also bounded by the right-hand side of (25).

To prove (26), we use the triangle inequality to infer that

$$\|\mathbf{K}_T^{-1/2}(\mathfrak{C}_T^{k+1} \underline{\sigma}_T - \sigma)\|_T \leq \|\mathbf{K}_T^{-1/2} \mathfrak{C}_T^{k+1}(\underline{\sigma}_T - \hat{\underline{\sigma}}_T)\|_T + \|\mathbf{K}_T^{-1/2}(\mathfrak{C}_T^{k+1} \hat{\underline{\sigma}}_T - \sigma)\|_T,$$

and bound the first term in the right-hand side using (25) since $\|\mathbf{K}_T^{-1/2} \mathfrak{C}_T^{k+1}(\underline{\sigma}_T - \hat{\underline{\sigma}}_T)\|_T \leq \|\underline{\sigma}_T - \hat{\underline{\sigma}}_T\|_{H,T}$ and the second term using (40).

Remark 10 (Potential error estimate). *Since $\sum_{T \in \mathcal{T}_h} (D_T^k \underline{\tau}_T, \hat{u}_T - u_T)_T = H(\underline{\sigma}_T - \hat{\underline{\sigma}}_T, \underline{\tau}_T) + \mathcal{E}_T(\underline{\tau}_T)$, a potential L^2 -error estimate of order $(k+1)$ bounding $\|\hat{u}_h - u_h\|$ classically follows from Lemma 5 and Theorem 6.*

4.4 Proof of Theorem 7

Let $z \in H_0^1(\Omega)$ be the unique solution of $-\nabla \cdot (\mathbf{K} \nabla z) = u_h - \hat{u}_h$, set $\omega := \mathbf{K} \nabla z$, and define $\hat{\omega}_h \in \underline{\Sigma}_h^k$, $\tilde{u}_h \in \mathbb{P}_d^{k+1}(\mathcal{T}_h)$, and $\tilde{z}_h \in \mathbb{P}_d^{k+1}(\mathcal{T}_h)$ such that, for all $T \in \mathcal{T}_h$,

$$\hat{\omega}_T := \underline{I}_T^k \omega, \quad \tilde{u}_T := \pi_T^{k+1} u, \quad \tilde{z}_T := \pi_T^{k+1} z.$$

We have, using $\nabla \cdot \omega = \hat{u}_h - u_h$ followed by $\nabla \cdot \omega \in \mathbb{P}_d^k(\mathcal{T}_h)$, $D_T^k \hat{\omega}_T = \nabla \cdot (\mathbf{K} \nabla z)|_T$ for all $T \in \mathcal{T}_h$ (a consequence of (12)), and (20a),

$$\begin{aligned} \|\hat{u}_h - u_h\|^2 &= (\hat{u}_h - u_h, \nabla \cdot \omega) = (u, \nabla \cdot \omega) - \sum_{T \in \mathcal{T}_h} (u_h, D_T^k \hat{\omega}_T) \\ &= \sum_{T \in \mathcal{T}_h} \{-(\sigma, \nabla z)_T + H_T(\underline{\sigma}_T, \hat{\omega}_T)\} := \mathfrak{T}_1 + \dots + \mathfrak{T}_5, \end{aligned}$$

with, for all $i \in \{1, \dots, 5\}$, $\mathfrak{T}_i = \sum_{T \in \mathcal{T}_h} \mathfrak{T}_{i,T}$ with

$$\begin{aligned} \mathfrak{T}_{1,T} &:= (\mathfrak{C}_T^{k+1} \underline{\sigma}_T - \sigma, \nabla(z - \tilde{z}_T))_T \\ \mathfrak{T}_{2,T} &:= (\mathbf{K}_T^{-1}(\mathfrak{C}_T^{k+1} \underline{\sigma}_T - \mathbf{K}_T \nabla \tilde{u}_T), \mathfrak{C}_T^{k+1} \hat{\omega}_T - \omega)_T, \\ \mathfrak{T}_{3,T} &:= S_T(\underline{\sigma}_T - \hat{\underline{\sigma}}_T, \hat{\omega}_T) + S_T(\hat{\underline{\sigma}}_T, \hat{\omega}_T), \\ \mathfrak{T}_{4,T} &:= (\mathfrak{C}_T^{k+1} \underline{\sigma}_T, \nabla \tilde{z}_T)_T - (\sigma, \nabla \tilde{z}_T)_T, \\ \mathfrak{T}_{5,T} &:= (\nabla \tilde{u}_T, \mathfrak{C}_T^{k+1} \hat{\omega}_T - \omega)_T. \end{aligned}$$

For the first term, the Cauchy–Schwarz inequality, the flux estimate (26), and the approximation properties (7) of π_T^{k+1} yield

$$|\mathfrak{T}_1| \lesssim \left\{ \sum_{T \in \mathcal{T}_h} \rho_{\mathbf{K},T} \lambda_{\sharp,T} h_T^{2(k+1)} |u|_{H^{k+2}(T)}^2 \right\}^{1/2} \lambda_{\sharp}^{1/2} h \|z\|_{H^2(\Omega)}.$$

For the second term, we use the Cauchy–Schwarz inequality to infer that

$$|\mathfrak{I}_2| \leq \left\{ \sum_{T \in \mathcal{T}_h} \|\mathbf{K}_T^{-1/2}(\mathfrak{C}_T^{k+1} \underline{\boldsymbol{\sigma}}_T - \mathbf{K}_T \nabla \check{u}_T)\|_T^2 \right\}^{1/2} \times \left\{ \sum_{T \in \mathcal{T}_h} \|\mathbf{K}_T^{-1/2}(\mathfrak{C}_T^{k+1} \hat{\boldsymbol{\omega}}_T - \boldsymbol{\omega})\|_T^2 \right\}^{1/2}.$$

To bound the first factor, we add and subtract $\boldsymbol{\sigma}$, use the triangle inequality, the flux estimate (25), and the approximation property (7) of π_T^{k+1} together with $\rho_{\mathbf{K},T} \geq 1$. To bound the second factor, we use the approximation property (40) of \mathfrak{C}_T^{k+1} . This yields

$$|\mathfrak{I}_2| \lesssim \left\{ \sum_{T \in \mathcal{T}_h} \rho_{\mathbf{K},T} \lambda_{\sharp,T} h_T^{2(k+1)} |u|_{H^{k+2}(T)}^2 \right\}^{1/2} \rho_{\mathbf{K}}^{1/2} \lambda_{\sharp}^{1/2} h \|z\|_{H^2(\Omega)}.$$

For the third term, the Cauchy–Schwarz inequality together with the flux estimate (25) and the approximation property (40) of S_T yield

$$|\mathfrak{I}_3| \lesssim \left\{ \sum_{T \in \mathcal{T}_h} \rho_{\mathbf{K},T} \lambda_{\sharp,T} h_T^{2(k+1)} |u|_{H^{k+2}(T)}^2 \right\}^{1/2} \rho_{\mathbf{K}}^{1/2} \lambda_{\sharp}^{1/2} h \|z\|_{H^2(\Omega)}.$$

For the fourth term, recalling (13) for the first summand, integrating by parts on T the second summand, and using the fact that $\nabla \cdot \boldsymbol{\sigma} = -f$ and $D_T^k \underline{\boldsymbol{\sigma}}_T = -\pi_T^k f$ owing to (2b) and (20b), respectively, we infer that

$$\mathfrak{I}_4 = \sum_{T \in \mathcal{T}_h} \left\{ (\pi_T^k f - f, \check{z}_T)_T + \sum_{F \in \mathcal{F}_T} (\sigma_F \epsilon_{TF} - \boldsymbol{\sigma} \cdot \mathbf{n}_{TF}, \check{z}_T)_F \right\}. \quad (47)$$

When $k = 0$, we estimate the first term in braces as follows:

$$|(\pi_T^k f - f, \check{z}_T)_T| = |(\pi_T^0 f - f, \check{z}_T - \pi_T^0 z)_T| \lesssim h_T^2 \|f\|_{H^1(T)} \|z\|_{H^1(T)},$$

whereas, for $k \geq 1$, we obtain

$$|(\pi_T^k f - f, \check{z}_T)_T| = |(\pi_T^k f - f, \check{z}_T - \pi_T^1 z)_T| \lesssim h_T^{k+2} \|f\|_{H^k(T)} \|z\|_{H^2(T)}.$$

Moreover, using the fact that σ_F , $\boldsymbol{\sigma} \cdot \mathbf{n}_F$, and z are single-valued at interfaces together with the fact that z vanishes on $\partial\Omega$, we can replace \check{z}_T by $(\check{z}_T - z)$ in the second term in braces in (47) to infer that

$$\left| \sum_{T \in \mathcal{T}_h} \sum_{F \in \mathcal{F}_T} (\sigma_F \epsilon_{TF} - \boldsymbol{\sigma} \cdot \mathbf{n}_{TF}, \check{z}_T - z)_F \right| \lesssim \lambda_{\sharp}^{1/2} \left\{ \sum_{T \in \mathcal{T}_h} \lambda_{\sharp,T} h_T^{2(k+2)} |u|_{H^{k+2}(T)}^2 \right\}^{1/2} h \|z\|_{H^2(\Omega)}.$$

where the last bound follows observing that $\|\boldsymbol{\sigma}\|_{H^{k+1}(\mathcal{T}_h)} \lesssim \lambda_{\sharp} |u|_{H^{k+2}(\mathcal{T}_h)}$. In conclusion, the following bound holds:

$$|\mathfrak{I}_4| \lesssim \left\{ \sum_{T \in \mathcal{T}_h} \lambda_{\sharp,T} h_T^{2(k+1)} |u|_{H^{k+2}(T)}^2 \right\}^{1/2} \lambda_{\sharp}^{1/2} h \|z\|_{H^2(\Omega)} + h^{k+2} \|f\|_{H^{k+\delta}(\mathcal{T}_h)} \|z\|_{H^2(\Omega)}.$$

For the fifth term, we proceed similarly. Using the definition of \mathfrak{C}_T^{k+1} , integration by parts, and the fact that $D_T^k \hat{\omega}_T = \nabla \cdot \omega = u_h - \hat{u}_h$ owing to the commuting property (12), we infer that

$$\mathfrak{T}_5 = \sum_{T \in \mathcal{T}_h} \sum_{F \in \mathcal{F}_T} (\check{u}_T, \hat{\omega}_F \epsilon_{TF} - \omega \cdot \mathbf{n}_{TF})_F = \sum_{T \in \mathcal{T}_h} \sum_{F \in \mathcal{F}_T} (\check{u}_T - u, \hat{\omega}_F \epsilon_{TF} - \omega \cdot \mathbf{n}_{TF})_F,$$

where we have used the fact that $\hat{\omega}_F$, $\omega \cdot \mathbf{n}_F$, and u are single-valued at interfaces together with the fact that u vanishes on $\partial\Omega$ to write $\check{u}_T - u$ in place of \check{u}_T . The Cauchy–Schwarz inequality and the approximation property (7) of π_T^{k+1} yield

$$|\mathfrak{T}_5| \lesssim \left\{ \sum_{T \in \mathcal{T}_h} \lambda_{\sharp, T} h_T^{2(k+1)} |u|_{H^{k+2}(T)}^2 \right\}^{1/2} \lambda_{\sharp}^{1/2} h \|z\|_{H^2(\Omega)}.$$

Gathering the above bounds for $\mathfrak{T}_1, \dots, \mathfrak{T}_5$, using elliptic regularity and recalling that $\rho_{\mathbf{K}, T} \geq 1$ for all $T \in \mathcal{T}_h$ and $\rho_{\mathbf{K}} \geq 1$, we obtain the estimate (27).

5 Numerical results

We present a numerical example for the homogeneous Dirichlet problem (2) on the unit square $\Omega = (0, 1)^2$ with diagonal diffusion tensor $\mathbf{K} = \begin{pmatrix} 1 & 0 \\ 0 & \rho_{\mathbf{K}}^{-1} \end{pmatrix}$ and exact solution $u = \sin(\pi x_1) \sin(\pi x_2)$.

We first evaluate the convergence rates for polynomial orders $0 \leq k \leq 4$ by solving the problem with $\rho_{\mathbf{K}} = 1$ on three mesh families obtained by refinement of the meshes depicted in Figure 1. Both the Kershaw and hexagonal meshes are obtained starting from a Cartesian grid. For the Kershaw mesh, only deformation is applied. For the hexagonal mesh, the connectivity is obtained starting from the nodes of a deformed Cartesian grid, and prescribing the hexagonal connectivity by selectively eliminating some of the nodes. In both cases, the mesh sequence is obtained by refining the underlying Cartesian grid and repeating the above procedure. The convergence results displayed in Figure 2 confirm the theoretical predictions of both Theorems 6 and 7.

We next evaluate numerically the dependence of the multiplicative constant in the error estimates on the anisotropy ratio $\rho_{\mathbf{K}}$ by solving the above problem on a fixed mesh with $\rho_{\mathbf{K}} \in \{2^i\}_{0 \leq i \leq 10}$. The mesh sizes are selected as follows: $7.68 \cdot 10^{-3}$ (triangular), $1.19 \cdot 10^{-2}$ (Kershaw), and $1.72 \cdot 10^{-2}$ (hexagonal). The results collected in Figure 3 show that the present method behaves in a somewhat more robust manner with respect to anisotropy than that predicted by the error estimates, in particular for the higher-orders and the hexagonal mesh family.

Acknowledgements. This work was partially supported by the ANR HHOMM project

References

- [1] P. F. Antonietti, S. Giani, and P. Houston. *hp*-version composite discontinuous Galerkin methods for elliptic problems on complicated domains. *SIAM J. Sci. Comput.*, 35(3):A1417–A1439, 2013.
- [2] R. Araya, C. Harder, D. Paredes, and F. Valentin. Multiscale hybrid-mixed method. *SIAM J. Numer. Anal.*, 51(6):3505–3531, 2013.

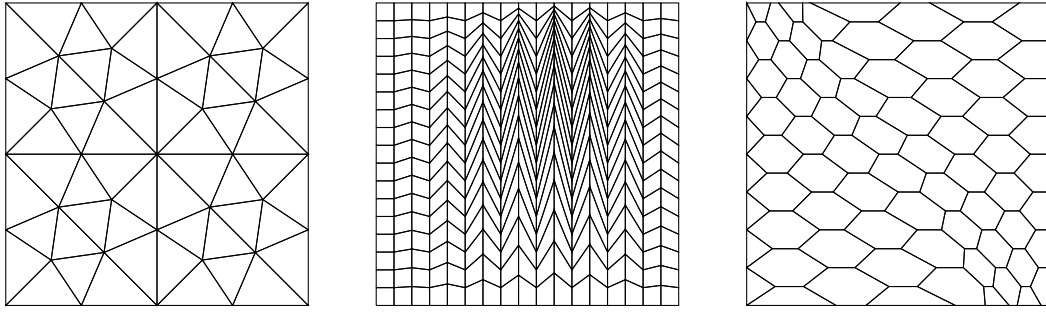
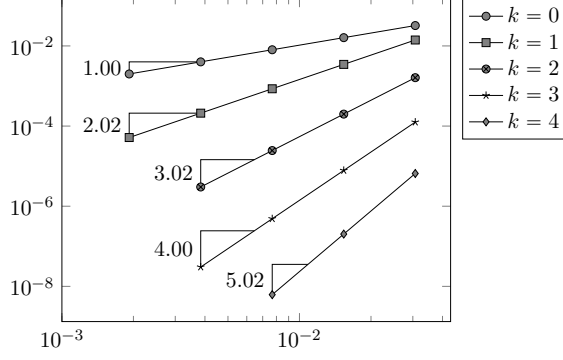
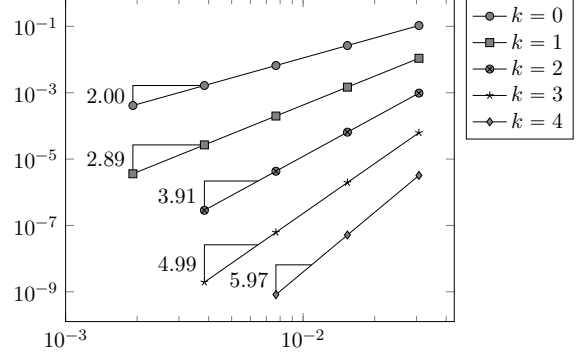


Figure 1: Triangular, Kershaw, and hexagonal meshes for the numerical example of Section 5

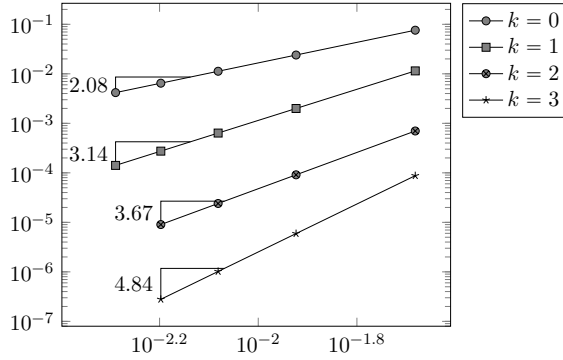
- [3] D. N. Arnold and F. Brezzi. Mixed and nonconforming finite element methods: implementation, postprocessing and error estimates. *RAIRO Modél. Math. Anal. Num.*, 19(4):7–32, 1985.
- [4] F. Bassi, L. Botti, A. Colombo, D. A. Di Pietro, and P. Tesini. On the flexibility of agglomeration based physical space discontinuous Galerkin discretizations. *J. Comput. Phys.*, 231(1):45–65, 2012.
- [5] L. Beirão da Veiga, F. Brezzi, A. Cangiani, G. Manzini, L. D. Marini, and A. Russo. Basic principles of virtual element methods. *Math. Models Methods Appl. Sci. (M3AS)*, 199(23):199–214, 2013.
- [6] L. Beirão da Veiga, F. Brezzi, and L. D. Marini. Virtual elements for linear elasticity problems. *SIAM J. Numer. Anal.*, 2(51):794–812, 2013.
- [7] L. Beirão da Veiga, F. Brezzi, L. D. Marini, and A. Russo. $H(\text{div})$ and $H(\text{curl})$ -conforming VEM. Technical report, IMATI-CNR Tech. Report 9PV14/0/0, 2014.
- [8] L. Beirão da Veiga, K. Lipnikov, and G. Manzini. Arbitrary-order nodal mimetic discretizations of elliptic problems on general meshes. *SIAM J. Numer. Anal.*, 5(49):1737–1760, 2011.
- [9] L. Beirão da Veiga, K. Lipnikov, and G. Manzini. *The Mimetic Finite Difference Method for Elliptic Problems*, volume 11 of *Modeling, Simulation and Applications*. Springer, 2014.
- [10] D. Boffi, F. Brezzi, and M. Fortin. *Mixed finite element methods and applications*, volume 44 of *Springer Series in Computational Mathematics*. Springer, Heidelberg, 2013.
- [11] J. Bonelle, D. A. Di Pietro, and A. Ern. Low-order reconstruction operators on polyhedral meshes: Application to Compatible Discrete Operator schemes. *Computer Aided Geometric Design*, 35–36:27–41, 2015.
- [12] J. Bonelle and A. Ern. Analysis of compatible discrete operator schemes for elliptic problems on polyhedral meshes. *ESAIM: Math. Model. Numer. Anal. (M2AN)*, 48:553–581, 2014.
- [13] J. Bonelle and A. Ern. Analysis of compatible discrete operator schemes for the Stokes equations on polyhedral meshes. *IMA J. Numer. Anal.*, 2015. DOI 10.1093/imanum/dru051.
- [14] A. Bossavit. A uniform rationale for Whitney forms on various supporting shapes. *Math. Comput. Simulation*, 80(8):1567–1577, 2010.
- [15] F. Brezzi, A. Buffa, and K. Lipnikov. Mimetic finite difference for elliptic problem. *ESAIM: Math. Model. Numer. Anal. (M2AN)*, 43:277–295, 2009.
- [16] F. Brezzi, A. Buffa, and G. Manzini. Mimetic scalar products of discrete differential forms. *J. Comput. Phys.*, 257:1228–1259, 2014.
- [17] F. Brezzi, R. S. Falk, and L. D. Marini. Basic principles of mixed virtual element methods. *ESAIM Math. Model. Numer. Anal.*, 48(4):1227–1240, 2014.



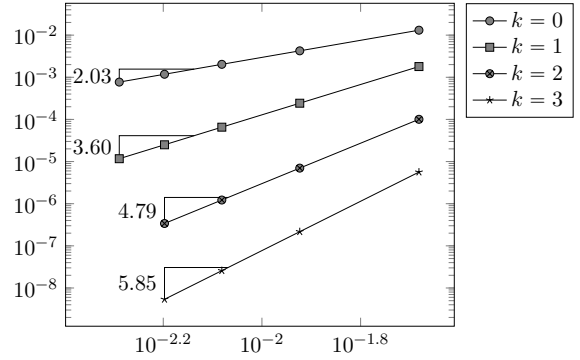
(a) $\|\hat{\sigma}_h - \underline{\sigma}_h\|_H$ vs. h , triangular mesh



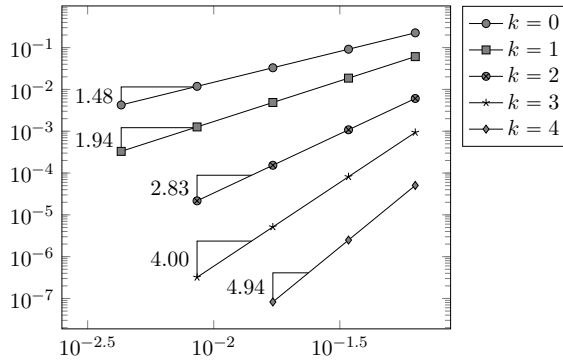
(b) $\|\hat{u}_h - u_h\|$ vs. h , triangular mesh



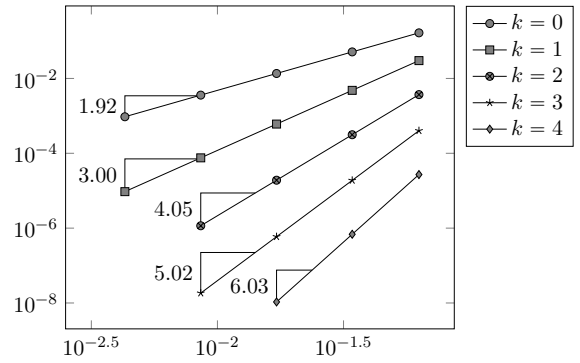
(c) $\|\hat{\sigma}_h - \underline{\sigma}_h\|_H$ vs. h , Kershaw mesh



(d) $\|\hat{u}_h - u_h\|$ vs. h , Kershaw mesh

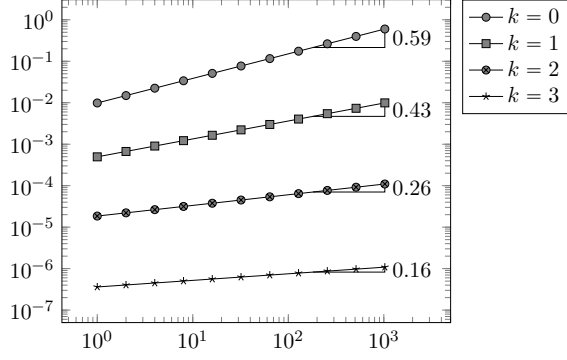


(e) $\|\hat{\sigma}_h - \underline{\sigma}_h\|_H$ vs. h , hexagonal mesh

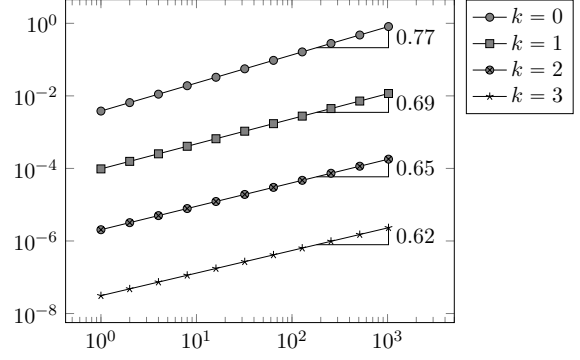


(f) $\|\hat{u}_h - u_h\|$ vs. h , hexagonal meshes

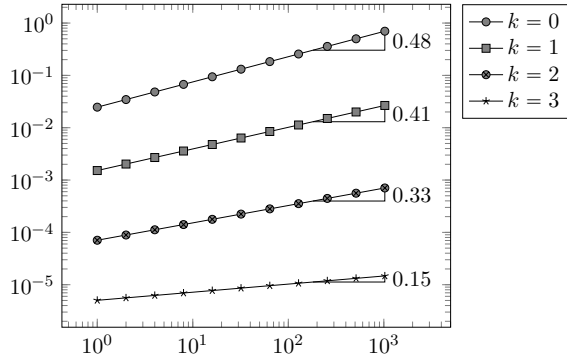
Figure 2: Convergence results for the numerical example of Section 5



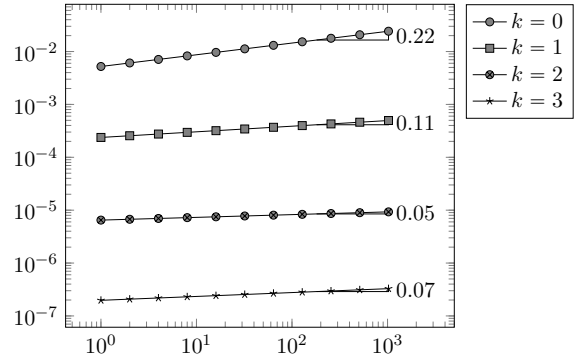
(a) $\|\hat{\sigma}_h - \underline{\sigma}_h\|_H$ vs. ρ_K , triangular mesh



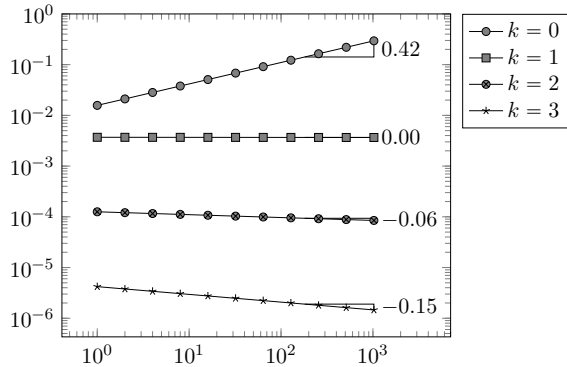
(b) $\|\hat{u}_h - u_h\|$ vs. ρ_K , triangular mesh



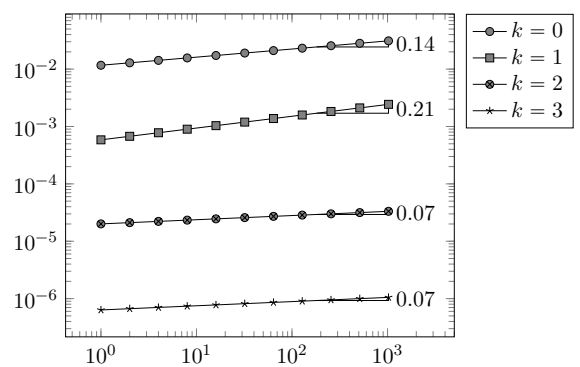
(c) $\|\hat{\sigma}_h - \underline{\sigma}_h\|_H$ vs. ρ_K , Kershaw mesh



(d) $\|\hat{u}_h - u_h\|$ vs. ρ_K , Kershaw mesh



(e) $\|\hat{\sigma}_h - \underline{\sigma}_h\|_H$ vs. ρ_K , hexagonal mesh



(f) $\|\hat{u}_h - u_h\|$ vs. ρ_K , hexagonal mesh

Figure 3: Error vs. ρ_K for a fixed mesh

- [18] F. Brezzi, K. Lipnikov, and M. Shashkov. Convergence of the mimetic finite difference method for diffusion problems on polyhedral meshes. *SIAM J. Numer. Anal.*, 43(5):1872–1896, 2005.
- [19] F. Brezzi, K. Lipnikov, M. Shashkov, and V. Simoncini. A new discretization methodology for diffusion problems on generalized polyhedral meshes. *Comput. Methods Appl. Mech. Engrg.*, 196(37–40):3682–3692, 2007.
- [20] A. Cangiani, E. H. Georgoulis, and P. Houston. *hp*-version discontinuous Galerkin methods on polygonal and polyhedral meshes. *Math. Models Methods Appl. Sci.*, 24(10):2009–2041, 2014.
- [21] S. H. Christiansen. A construction of spaces of compatible differential forms on cellular complexes. *Math. Models Methods Appl. Sci.*, 18(5):739–757, 2008.
- [22] B. Cockburn, D. A. Di Pietro, and A. Ern. Bridging the hybrid high-order and hybridizable discontinuous galerkin methods. *ESAIM: Math. Model. Numer. Anal. (M2AN)*, 2015. Accepted for publication.
- [23] B. Cockburn, J. Gopalakrishnan, and R. Lazarov. Unified hybridization of discontinuous Galerkin, mixed, and continuous Galerkin methods for second order elliptic problems. *SIAM J. Numer. Anal.*, 47(2):1319–1365, 2009.
- [24] B. Cockburn, J. Gopalakrishnan, and F.-J. Sayas. A projection-based error analysis of HDG methods. *Math. Comp.*, 79:1351–1367, 2010.
- [25] B. Cockburn, W. Qiu, and K. Shi. Conditions for superconvergence of HDG methods for second-order elliptic problems. *Math. Comp.*, 81:1327–1353, 2012.
- [26] L. Codecasa, R. Specogna, and F. Trevisan. A new set of basis functions for the discrete geometric approach. *J. Comput. Phys.*, 19(299):7401–7410, 2010.
- [27] D. A. Di Pietro. Cell centered Galerkin methods for diffusive problems. *ESAIM: Math. Model. Numer. Anal. (M2AN)*, 46(1):111–144, 2012.
- [28] D. A. Di Pietro and A. Ern. *Mathematical aspects of discontinuous Galerkin methods*, volume 69 of *Mathématiques & Applications*. Springer-Verlag, Berlin, 2012.
- [29] D. A. Di Pietro and A. Ern. A hybrid high-order locking-free method for linear elasticity on general meshes. *Comput. Meth. Appl. Mech. Engrg.*, 283:1–21, 2015.
- [30] D. A. Di Pietro and A. Ern. Hybrid high-order methods for variable-diffusion problems on general meshes. *C. R. Acad. Sci Paris, Ser. I*, 353:31–34, 2015.
- [31] D. A. Di Pietro, A. Ern, and S. Lemaire. An arbitrary-order and compact-stencil discretization of diffusion on general meshes based on local reconstruction operators. *Comput. Meth. Appl. Math.*, 14(4):461–472, 2014.
- [32] D. A. Di Pietro and S. Lemaire. An extension of the Crouzeix–Raviart space to general meshes with application to quasi-incompressible linear elasticity and Stokes flow. *Math. Comp.*, 84(291):1–31, 2015.
- [33] J. Droniou. Finite volume schemes for diffusion equations: introduction to and review of modern methods. *Math. Models Methods Appl. Sci.*, 24(8):1575–1619, 2014.
- [34] J. Droniou and R. Eymard. A mixed finite volume scheme for anisotropic diffusion problems on any grid. *Numer. Math.*, 105:35–71, 2006.
- [35] J. Droniou, R. Eymard, T. Gallouët, and R. Herbin. A unified approach to mimetic finite difference, hybrid finite volume and mixed finite volume methods. *Math. Models Methods Appl. Sci. (M3AS)*, 20(2):1–31, 2010.
- [36] J. Droniou, R. Eymard, T. Gallouët, and R. Herbin. Gradient schemes: a generic framework for the discretisation of linear, nonlinear and nonlocal elliptic and parabolic equations. *Math. Models Methods Appl. Sci.*, 23(13):2395–2432, 2013.

- [37] R. Eymard, T. Gallouët, and R. Herbin. Discretization of heterogeneous and anisotropic diffusion problems on general nonconforming meshes. SUSI: a scheme using stabilization and hybrid interfaces. *IMA J. Numer. Anal.*, 30(4):1009–1043, 2010.
- [38] R. Eymard, C. Guichard, and R. Herbin. Small-stencil 3D schemes for diffusive flows in porous media. *ESAIM Math. Model. Numer. Anal.*, 46(2):265–290, 2012.
- [39] C. Harder, D. Paredes, and F. Valentin. A family of multiscale hybrid-mixed finite element methods for the Darcy equation with rough coefficients. *J. Comput. Phys.*, 245:107–130, 2013.
- [40] Y. Kuznetsov, K. Lipnikov, and M. Shashkov. Mimetic finite difference method on polygonal meshes for diffusion-type problems. *Comput. Geosci.*, 8:301–324, 2004.
- [41] K. Lipnikov and G. Manzini. A high-order mimetic method on unstructured polyhedral meshes for the diffusion equation. *J. Comput. Phys.*, 272:360–385, 2014.
- [42] M. Vohralík and B. I. Wohlmuth. Mixed finite element methods: implementation with one unknown per element, local flux expressions, positivity, polygonal meshes, and relations to other methods. *Math. Models Methods Appl. Sci.*, 23(5):803–838, 2013.
- [43] J. Wang and X. Ye. A weak Galerkin finite element method for second-order elliptic problems. *J. Comput. Appl. Math.*, 241:103–115, 2013.
- [44] J. Wang and X. Ye. A weak Galerkin mixed finite element method for second order elliptic problems. *Math. Comp.*, 83(289):2101–2126, 2014.